

MCCMF: Collaborative matrix factorization based on matrix completion for predicting miRNA-disease associations

Tian-Ru Wu

Qufu Normal University - Rizhao Campus

Meng-Meng Yin

Qufu Normal University - Rizhao Campus

Cui-Na Jiao

Qufu Normal University - Rizhao Campus

Ying-Lian Gao

Qufu Normal University - Rizhao Campus

Xiang-Zhen Kong

Qufu Normal University - Rizhao Campus

Jin-Xing Liu (✉ sd cavell@126.com)

Qufu Normal University - Rizhao Campus <https://orcid.org/0000-0001-6104-2149>

Methodology article

Keywords: MiRNA-disease association prediction, Matrix completion, Weight K nearest known neighbors, Matrix factorization

Posted Date: October 13th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-36602/v4>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on October 14th, 2020. See the published version at <https://doi.org/10.1186/s12859-020-03799-6>.

MCCMF: Collaborative matrix factorization based on matrix completion for predicting miRNA-disease associations

**Tian-Ru Wu, Meng-Meng Yin, Cui-Na Jiao, Ying-Lian Gao, Xiang-Zhen Kong,
Jin-Xing Liu***

School of Computer Science, Qufu Normal University, Rizhao, 276826, China.

* Corresponding author

Email addresses:

TRW: wutianru@126.com

MMY: yinmengmeng00@163.com

CNJ: jiaocuina123@163.com

YLG: yinliangao@126.com

XZK: kongxzhen@163.com

JXL: sdcavell@126.com

Abstract

Background: microRNAs (miRNAs) are non-coding RNAs with regulatory functions. Many studies have shown that miRNAs are closely associated with human diseases. Among the methods to explore the relationship between the miRNA and the disease, traditional methods are time-consuming and the accuracy needs to be improved. In view of the shortcoming of previous models, a method, collaborative matrix factorization based on matrix completion (MCCMF) is proposed to predict the unknown miRNA-disease associations.

Results: The complete matrix of the miRNA and the disease is obtained by matrix completion. Moreover, Gaussian Interaction Profile (GIP) kernel is added to the miRNA functional similarity matrix and the disease semantic similarity matrix. Then the Weight K Nearest Known Neighbors (WKNKN) method is used to pretreat the association matrix, so the model is close to the reality. Finally, collaborative matrix factorization (CMF) method is applied to obtain the prediction results. Therefore, the MCCMF obtains a satisfactory result in the five-fold cross-validation, with an AUC of 0.9569(0.0005).

Conclusions: The AUC value of MCCMF is higher than other advanced methods in the 5-fold cross validation experiment. In order to comprehensively evaluate the performance of MCCMF, accuracy, precision, recall and f-measure are also added. The final experimental results demonstrate that MCCMF outperforms other methods in predicting miRNA-disease associations. In the end, the effectiveness and practicability of MCCMF are further verified by researching three specific diseases.

Keywords

MiRNA-disease association prediction, Matrix completion, Weight K nearest known neighbors, Matrix factorization

Background

MicroRNAs (MiRNAs) are a class of non-coding single-stranded RNA molecules. Their lengths are usually 18 to 24 nucleotides. Instead of synthesizing proteins, miRNAs participate in post-transcriptional regulation of gene expression in eukaryotes and viruses [1]. In spite of the first miRNA Line-4 was discovered in 1993 [2], the diversity and prevalence of these genes were revealed in recent years. To date, 38,589 miRNA have been found in animals, plants and viruses [3]. At the same time, miRNAs were discovered to play an important role in cell proliferation [4], differentiation [5], senescence [6], apoptosis [7], and so on. A study indicated that more than one third of human genes are regulated by miRNA [8]. Obviously, miRNA disorder could have severe impacts on humans.

Evidence shows that an increasing number of miRNAs are closely associated with diseases [9]. Since the first discovery of miR15 and miR16 deficiency in B cell chronic lymphocytic leukemia (B-CLL) [10], the research results of miRNA-disease associations are often reported. For example, the expression of miR-25 and miR-223 is significantly higher in patients with esophageal squamous cell carcinoma than the normal people, while the expression of miR-375 is significantly lower [11]. Studies show that miR-26a may be a regulatory factor that inhibits the progression and metastasis of c-Myc/EZH2 double height advanced HCC [12]. In addition, miR-340

has been suggested as a biomarker for cancer metastasis and prognosis [13]. At present, the research on miRNAs and diseases is becoming more extensive. Researchers have also developed a number of databases to store miRNA and disease data, such as dbDEMC [14], HMDD v3.0 [15] and miR2Disease [16]. Unfortunately, the known correlation data is not complete. Moreover, traditional methods to identify new miRNA-disease associations are time-consuming and laborious.

With the improvement of information technology and the development of a large number of miRNA data sets, many effective methods for predicting miRNA-disease associations have been proposed [17]. According to the hypothesis that functionally similar miRNAs may be associated with diseases with similar phenotypes [18], Jiang *et al.* first constructed a genetic data network, and then prioritized disease-related miRNAs to predict miRNA-disease associations [19]. However, due to the limited association information, this method is not quite effective. A computational framework was developed by Li *et al.*, which can be used to measure the association between the cancer and miRNA based on the functional consistency score (FCS) of miRNA target genes and cancer-related genes. This method has a significant advantage in the identification of cancer-related miRNA [20]. Based on heterogeneous omics data, the potential miRNA-disease associations were identified via using a Graph Regularized Non-negative Matrix Factorization (GRNMF) by Xiao *et al.* [21]. However, the prediction results of GRNMF method may not be optimal in some cases. Chen *et al.* proposed a new a computational model of Matrix Decomposition and Heterogeneous Graph Inference for miRNA-disease association prediction (MDHGI)

to discover new miRNA-disease associations [22]. The model made full use of matrix decomposition before the construction of heterogeneous networks, thus improving the prediction accuracy. The protein-driven inference of miRNA-disease associations (miRPD) was proposed by Mørk *et al.*, which can infer the correlation between miRNA-protein-disease associations [23]. At the same time, they provide scoring schemes that can create correlation sets of high and medium credibility. Three new miRNA-disease association prediction methods based on global network similarity measure were developed by Chen *et al.*, namely MBSI (microRNA-based similarity inference), PBSI (phenotype-based similarity inference) and NetCBI (network-consistency-based inference) [24]. NetCBI is especially suitable for predicting target diseases, but it relies on network similarity measurement to a great extent. Similarly, Gao *et al.* put forward a method based on Double Network Sparse Graph Regularized Matrix Factorization (DNSGRMF), and added the $L_{2,1}$ -norm and Gaussian interaction profile (GIP) kernel to improve the prediction ability [25]. In addition, considering the nearest neighbor information of the miRNA and the disease, Gao *et al.* introduced a method of Nearest Profile-based Collaborative Matrix Factorization (NPCMF) to predict miRNA-disease associations [26]. One of the most obvious disadvantages of NPCMF is that it introduces too much NP information, which may reduce the prediction accuracy while adding extra noise. In order to protect the known correlation, Logistic Weighted Profile-based Collaborative Matrix Factorization (LWPCMF) method was proposed by Yin *et al.*, which effectively predicts miRNA-disease associations [27]. The prediction effect of this method is

promising. Chen *et al.* constructed a model based on Canonical Correlation analysis (CCA), which can fully reveal the possible molecular causes of miRNA-disease association [28]. However, direct performance comparison is difficult to be achieved by this method.

In recent years, machine learning-based miRNA-disease association prediction methods are also popular. A support vector machine (SVM) classifier was developed by Xu *et al.* to extract features from the miRNA-disease network and miRNA expression levels [29]. Yet, the construction of miRNA target-dysregulated network (MTDN) is complex, so only direct miRNA target regulation can be predicted. Chen *et al.* used random walk to prioritize disease-related miRNAs to predict potential human miRNA-disease associations [30]. Like the problem of Jiang *et al.*, their approach is also affected by limited disease-miRNA associations. A model of Restricted Boltzmann machine for multiple types of miRNA-disease association prediction (RBMMMDA) was established by Chen *et al.*[31]. Chen *et al.* constructed a computational model called Laplacian Regularized Sparse Subspace Learning for MiRNA-Disease Association prediction (LRSSLMDA) [32]. The model has stronger dimensionality reduction capability and can be easily extended to higher dimensional data sets. A new Induction Matrix Completion model for MiRNA-Disease Association prediction (IMCMDA) was constructed by Chen *et al.* [33]. Because it is a semi-supervised model, only positive samples and unmarked samples are needed, which greatly reduces the difficulty of modeling. Soon after, Chen *et al.* proposed a new MiRNA-Disease Association Prediction Bipartite graph Network Projection

computing model (BNPMDA) [34]. Compared with previous models, the prediction accuracy of BNPMDA is improved. A new miRNA-disease association prediction algorithm based on the decision tree was proposed by Chen *et al.* [35]. This method constructs a computing framework for integrated learning and dimension reduction. By training and integrating multiple base classifiers, they reduce prediction bias and improve prediction performance. Ding *et al.* used an improved calculation method based on inductive matrix completion to predict miRNA-disease associations. (IIMCMP) [36]. Experiments show that IIMCMP can achieve powerful and reliable performance evaluation. Li *et al.* developed a method of neural inductive matrix completion with graph convolutional network (NIMCGCN) for the prediction of miRNA-disease association [37]. To test the predictive power of NIMCGCN in the absence of any known miRNAs, they studied breast cancer with 100% accuracy. The above methods have made great contributions to predicting associations of miRNA-disease.

Since the shortcomings of the above methods, a novel method for predicting miRNA-disease associations with Collaborative Matrix Factorization based on Matrix Completion (MCCMF) is proposed in this paper. Firstly, human miRNA-disease association matrix, miRNA function similarity matrix and disease semantic similarity matrix are obtained from HMDD v2.0, but the obtained matrix is sparse. Therefore, the matrix completion method is used to complete the matrix. The matrix completion algorithm is mainly developed on the basis of Augmented Lagrange multiplier method (ALM) [38], Alternating Direction Method (ADM) [39] and Singular Value Threshold

(SVT) operation [40]. Secondly, we integrate the completed matrix and the GIP kernel similarity matrix of the disease and the miRNA. At the same time, the miRNA-disease association matrix is preprocessed by Weight K Nearest Known Neighbors (WKNKN) method to solve the problem of unknown missing values [41]. Finally, collaborative matrix factorization is used to predict associations between miRNAs and diseases. In the experiment, a five-fold cross validation on MCCMF is performed, and results show that our method is superior to the other four methods. In addition, we focus on the cases of Gastrointestinal Neoplasms, Retinoblastoma and Hepatoblastoma. Our method not only successfully verifies the known associations of miRNA-disease, but also finds many unknown associations. To sum up, MCCMF can avoid the inherent noise of the data set, with high-speed and high prediction accuracy.

Results

Performance evaluation

In this section, AUC value, accuracy, precision, recall and f-measure are used to evaluate the performance of MCCMF method. Initially, we implement five-fold cross validation to objectively evaluate the predictive power of our method. The existing miRNA-disease associations are randomly divided into five groups, among which four groups are used as the training set and the remaining one as the test set. In addition, in order to demonstrate the high predictive capability of our method, the random deletion of the miRNA-disease association (i.e. Cross Validation pairs' mode) increases the difficulty of prediction before performing the cross validation [42]. Five-fold cross validation is repeated 10 times to prevent grouping from causing bias, and the average

result of 10 times is used as the final evaluation result.

The ROC curve is drawn to represent the predicted performance intuitively, and the AUC value is calculated to evaluate MCCMF quantitatively. TPR and FPR can be expressed as:

$$TPR = \frac{TP}{TP + FN}, \quad (1)$$

$$FPR = \frac{FP}{TN + FP}, \quad (2)$$

where TP is the number of samples that are actually positive and are also predicted to be positive. FN represents the number of samples that are actually negative and also predicted to be negative. However, TN and FP represent the number of samples for which the predicted results are inconsistent.

In order to make the performance evaluation more comprehensive, we also use other evaluation indicators, including the accuracy, precision, recall and f-measure. Their calculation formulas are defined as follows:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}, \quad (3)$$

$$precision = \frac{TP}{TP + FP}, \quad (4)$$

$$recall = \frac{TP}{TP + FN}, \quad (5)$$

$$f - measure = \frac{2 \times precision \times recall}{precision + recall}. \quad (6)$$

Comparison with Other Methods

The AUC value is generally between 0 and 1. The higher the AUC value is, the better the prediction result will be. MCCMF finally obtains an AUC value of 0.9563 in the five-fold cross validation. MCCMF is compared with four advanced methods such as WBNPMD [43], RLSMDA [44], GRNMF [21] and CMF [45], which proves the

superior performance of our method. The ROC curves are drawn in Fig.1, and the comparison results are listed in Table 1. The results of other methods in Table 1 are obtained directly from the literature.

In the Table 1, the highest value is highlighted in bold, with the standard deviation in parentheses. In the five-fold cross validation experiment, WBNPMD, RLSMDA, GRNMF, CMF and MCCMF obtain AUCs of 0.9173, 0.8389, 0.869, 0.8697 and 0.9569, respectively. Therefore, our method is superior to the other four methods.

WBNPMD with higher AUC value is selected for comparison with MCCMF, and accuracy, precision, recall and f-measure are presented as a bar graph in Fig.2. Also, MCCMF is better than WBNPMD.

Case Studies

In the end, we carry out a simulation experiment to analyze the specific disease. First of all, the disease we want to explore is selected and the predicted score is ranked. Then, based on the predicted score after ranking, some miRNAs of high associations degree with the disease are found. Moreover, by comparing with the original miRNA-disease association matrix, they are determined whether the associations of high prediction score is known. Finally, the unknown associations are verified by searching existing data sets. Here, we choose three diseases of Gastrointestinal Neoplasms, Retinoblastoma and Hepatoblastoma for analysis. In addition, three popular data sets, dbDEMC [14], HMDD v3.0 [46] and miRCancer [47] are used for validation. These data sets store miRNA-disease associations that have been experimentally confirmed by some researchers over the years.

Gastrointestinal Neoplasms is a very common gastrointestinal disease with a high

incidence. However, there are no obvious symptoms in the early growth stage of the neoplasms, which is very dangerous to human beings. We successfully predict 31 known associations and 9 new associations, 7 of which are confirmed by HMDD v3.0 and miRCancer. For example, Tazawa *et al.* discovered the potential role of oncogenic miR-21 in Gastrointestinal Neoplasms [48]. Other confirmed miRNAs have been reported in relevant data sets, and they are not listed here. There are still two unconfirmed ones that need further research. Table 2 describes the simulation results, where known associations are shown in bold, confirmed new predictions are written to the corresponding database, and unconfirmed ones are shown as “unconfirmed”. The predicted scores in the Table 2 are ranked according to the strength of the association between the miRNA and disease. There is a threshold to determine whether the prediction is accurate. Compared with known information and other databases, the prediction results of our method are generally accurate. Although two remain unconfirmed, these two could provide some insights for researchers.

Retinoblastoma is a malignant tumor that occurs in children under 3 years old, and has a familial predisposition. There are 38 known associations between the disease and miRNA in the known association data set, and 37 known associations are successfully predicted by us. At the same time, 23 new associations are predicted, seven of which are confirmed and the others are unconfirmed. Montoya *et al.* found that the expression of miR-31 in Retinoblastoma is significantly reduced, which promotes the development of targeted therapy for Retinoblastoma [49]. Table 3 shows the specific situation. The predictive sorting method in Table 3 is the same as that in Table 2.

Hepatoblastoma is the most common intraabdominal malignant tumor after neuroblastoma and nephroblastoma in childhood. In the existing miRNA-disease association data set, there are 8 known miRNA-disease associations, and all of them have been predicted. Besides, we predicted 12 new associations, seven of which are confirmed and 5 are not. We also find literatures confirming that miR-143 is a factor affecting Hepatoblastoma. The study of Zhang *et al.* showed that blocking miR-143 could significantly inhibit local liver metastasis [50]. Hepatoblastoma prediction results are shown in Table 4. The predictive sorting method in Table 4 is also the same as that in Table 2 and 3.

As can be seen from the simulation results above, most known miRNAs are successfully predicted, while a small number of unknown associations are in HMDDv3.0, miRCancer and dbDEMC data sets. Although a few have not been confirmed, they can be used as a reference for researchers. In addition, we used Cytoscape software to map the prediction network of these three diseases (Fig.3). In the network, the ellipse represents miRNAs, and the remaining shapes represent diseases. The correlations are connected by line segments with arrows, and there are common miRNAs between diseases. According to the size of the predicted score, the color degree of the ellipse is set differently. The darker the color of the ellipse is set to, the stronger the correlation between miRNA and disease is.

Discussion

The above experimental results are enough to prove that our method is superior to the most advanced method. The excellent prediction performance of MCCMF can be

attributed to several significant factors. Firstly, data is preprocessed by Weight K Nearest Known Neighbors method and matrix completion method to improve the prediction accuracy. Secondly, a collaborative matrix factorization model is applied to predicting miRNA-disease associations, which is a promising one among many collaborative filtering technologies. In bioinformatics, matrix factorization contributes to identifying hidden links among genes. However, the performance of our method needs to be further improved. For instance, there exists a better way to integrate data, rather than simply adding them together. In the future, we will improve the technology to use the latest version of the data set, such as HMDD v3.0.

Conclusions

In this paper, a collaborative matrix factorization method based on matrix completion (MCCMF) is developed for predicting miRNA-disease associations. Considering the sparse and incomplete similarity matrix of miRNA-disease, we use the matrix completion method to complete the matrix. Then the completed matrix is integrated with GIP kernel similarity to improve the data information and reduce the influence of noises. In addition, WKNKN is also introduced to pretreat the existing association matrix of miRNAs and diseases, so our method is suitable to practical problems. Finally, the idea of CMF is adopted to construct the objective function and obtain the predicted results. The AUC value (0.9569) of MCCMF is higher than other advanced methods in the 5-fold cross validation experiment. In order to comprehensively evaluate the performance of MCCMF, accuracy, precision, recall and f-measure are applied to measure the performance, and results are 0.992, 0.779, 0.918 and 0.830,

respectively. Compared with the other four methods, our method has the best performance. The analysis of Gastrointestinal Neoplasms, Retinoblastoma and Hepatoblastoma further verified the effectiveness of MCCMF. Since most of associations are unknown in reality, MCCMF can also be used to predict in this situation.

Methods

We develop a novel method for predicting miRNA-disease associations with MCCMF. MCCMF is divided into four main steps: Firstly, we use the matrix completion algorithm to complete the miRNA similarity matrix and the disease similarity matrix to generate a new completion similarity matrix. Secondly, the new completion similarity matrix is integrated with existing miRNA and disease similarity information. Thirdly, the WKNKN is used to convert the binary values of the miRNA-disease association matrix into the interaction likelihood values [41]. Finally, the Collaborative Matrix Factorization is used to predict the association of miRNA-disease. Fig. 4 shows the complete process for MCCMF.

Human miRNA-disease associations

The initial miRNA-disease association data is downloaded from HMDD v2.0 [51]. HMDD v2.0 is an experimental data set supporting human miRNA-disease associations, and storing 5430 experimentally verified miRNA-disease associations between 495 miRNAs and 383 diseases. In this paper, the adjacency matrix \mathbf{MD} is used to represent the miRNA-disease association network. The adjacency matrix \mathbf{MD} is a sparse matrix composed of 0 and 1. If $\mathbf{MD}(m_i, d_j)$ is 1, disease d_j is correlated

with miRNA m_i ; otherwise irrelevant.

MiRNA function similarity

According to the hypothesis that functionally similar miRNAs are more likely to be associated with phenotypic diseases, a method for calculating the functional similarity of miRNAs (MISIM) is proposed by Wang *et al.* [52]. Firstly, we need to define semantic similarity between one disease and one group of disease. The calculation formula is as follows:

$$S(d, \mathbf{D}) = \max_{1 \leq i \leq k} (S(d, \mathbf{D}_i)). \quad (7)$$

Here d represents one disease and \mathbf{D} represents one disease group. Then, we define the similarity of d and \mathbf{D} , $S(d, \mathbf{D})$, as the maximum similarity.

Functional similarity of the two miRNAs is defined as

$$MISIM(M_1, M_2) = \frac{\sum_{1 \leq i \leq m} S(d_{1i}, \mathbf{D}_2) + \sum_{1 \leq j \leq n} S(d_{2j}, \mathbf{D}_1)}{m + n}, \quad (8)$$

where M_1 and M_2 represent the related miRNAs of \mathbf{D}_1 and \mathbf{D}_2 , respectively.

\mathbf{D}_1 contains m diseases, and \mathbf{D}_2 contains n diseases.

In this paper, we download the miRNA function similarity from <http://www.cuilab.cn/files/images/cuilab/misim.zip>. And the matrix \mathbf{MF} is used to represent the functional similarity network of the miRNA, in which the element $\mathbf{MF}(i, j)$ represents the similarity between miRNA m_i and miRNA m_j . The self-similarity of each miRNA is 1, so the diagonal elements of the matrix \mathbf{MF} are 1.

Due to incomplete miRNA data supported by the experiment, the similarity values calculated by MISIM may be biased. Some subsequent treatment of the matrix may be improved [53].

Disease semantic similarity

The relationship between different diseases is obtained from the MeSH database (<http://www.ncbi.nlm.nih.gov/>). Based on the previous literature [52], we represent the disease D as a Directed Acyclic Graph, $DAG(D) = (D, T(D), E(D))$, where $T(D)$ is the set of both a node D and its ancestor nodes, and $E(D)$ is the set of edges that ancestor nodes pointing to node D . For ancestor node t in $DAG(A)$, its contribution to the semantic value of disease A is computed as follows:

$$Dl_A(t) = \begin{cases} 1 & \text{if } t = A, \\ \max \{ \Delta * Dl_A(t') \mid t' \in \text{children of } t \} & \text{if } t \neq A. \end{cases} \quad (9)$$

In the above formula, Δ is a semantic contribution factor. Based on the method of Wang *et al.*, the value of Δ is set to 0.5. For the disease A , the contribution of itself to the disease A is 1, while the contribution of ancestor node t is decreasing with the increase of its layers.

Based on the contribution of ancestor diseases and disease A itself, the semantic value of disease A can be expressed as follows:

$$DV1(A) = \sum_{t \in T(A)} Dl_A(t). \quad (10)$$

According to the hypothesis that the more shared part of the disease pairs in $DAGs$ is, the higher similarity is. The semantic similarity between disease A and disease B is calculated as:

$$DS1(A, B) = \frac{\sum_{t \in T(A) \cap T(B)} (Dl_A(t) + Dl_B(t))}{DV1(A) + DV1(B)}. \quad (11)$$

However, the above model is a little inadequacy, which is the setting of Δ that causes the same layer of diseases with the same semantic contribution. Obviously, the incidence of various diseases is different, and the contribution of diseases with high

incidence should be less than those with low incidence. To improve the above model, we combine the method of Xuan *et al.* to define the semantic similarity calculation method [54]. In this method, the contribution of ancestor node t in $DAG(A)$ to the semantic value of disease A is as follows:

$$D2_A(t) = -\log \frac{\text{the number of DAGs including } t}{\text{the number of diseases}}. \quad (12)$$

The semantic value of disease A , and the semantic similarity between the disease A and the disease B are calculated as:

$$DV2(A) = \sum_{t \in T(A)} D2_A(t), \quad (13)$$

$$DS2(A, B) = \frac{\sum_{t \in T(A) \cap T(B)} (D2_A(t) + D2_B(t))}{DV2(A) + DV2(B)}. \quad (14)$$

Finally, in order to calculate the semantic similarity more comprehensive and rational, we combine the two models to get Equation (15).

$$DS(A, B) = \frac{DS1(A, B) + DS2(A, B)}{2}. \quad (15)$$

Gaussian interaction profile kernel similarity for diseases and miRNAs

On the basis of the hypothesis that functionally similar miRNAs may be associated with similar diseases, and vice versa, the known miRNA-disease association network is used to construct the GIP kernel similarity for diseases and miRNAs [55]. GIP kernel similarity can increase the multiple and topological information of known correlations. The interaction profile of miRNA $m(i)$ is represented by the binary vector $M(i)$ of the i -th column of the adjacency matrix \mathbf{MD} . Similarly, the binary vector $D(i)$ of the i -th row of the adjacency matrix \mathbf{MD} denotes the interaction profile of disease $d(i)$. Hence, we can define the GIP kernel similarity for miRNAs and diseases as

follows:

$$\mathbf{GM}(m(i), m(j)) = \exp(-\gamma_m \|M(i) - M(j)\|^2), \quad (16)$$

$$\mathbf{GD}(d(i), d(j)) = \exp(-\gamma_d \|D(i) - D(j)\|^2). \quad (17)$$

Here, γ_m and γ_d are parameters to control the kernel bandwidth and obtained by the following formulas:

$$\gamma_m = \frac{\delta_m}{\frac{1}{nm} \sum_{i=1}^{nm} \|M(i)\|^2}, \quad (18)$$

$$\gamma_d = \frac{\delta_d}{\frac{1}{nd} \sum_{i=1}^{nd} \|D(i)\|^2}, \quad (19)$$

where δ_m and δ_d are also bandwidth parameters and they are set to 1 according to the previous study [56]. The nm and nd mean the number of all the miRNAs and diseases.

Matrix completion

The miRNA functional similarity matrix and disease semantic similarity matrix calculated by the above operations are still sparse and incomplete, and there are some redundant associations (i.e. inherent noise). So we use the matrix completion method to solve the problem [57]. Suppose the incomplete matrix is \mathbf{D} , which can be represented as a linear combination of \mathbf{D} and the noise matrix \mathbf{N} . The formula is as follows:

$$\mathbf{D} = \mathbf{DR} + \mathbf{N}, \quad (20)$$

where \mathbf{DR} is a low-rank matrix, and specifically, it is a more refined or informative similarity matrix after removing noise from the existing similarity matrix.

In order to make \mathbf{R} be low-rank, a nuclear norm on \mathbf{D} is added. At the same time, the $L_{2,1}$ -norm of the error term \mathbf{N} is used to make noise matrix \mathbf{N} more sparse.

When the final low-rank matrix \mathbf{DR}^* and sparse matrix \mathbf{N}^* are calculated, \mathbf{DR}^* or $\mathbf{D-N}^*$ are used to describe a completed matrix. Therefore, a formula for solving convex optimization problem can be defined as follows:

$$\min_{\mathbf{R}, \mathbf{N}} \|\mathbf{R}\|_* + \omega \|\mathbf{N}\|_{2,1} \text{ s.t. } \mathbf{D} = \mathbf{DR} + \mathbf{N}. \quad (21)$$

Here, $\|\cdot\|_*$ represents the nuclear norm, $\omega \in (0,1)$ is the positive weighting parameter and $\|\cdot\|_{2,1}$ is the noise regularization term.

When solving optimization problems under equality constraints, the ALM method is more effective [38]. Therefore, according to ALM, the Equation (21) can be rewritten as:

$$\min_{\mathbf{R}, \mathbf{N}, \mathbf{X}} \|\mathbf{X}\|_* + \omega \|\mathbf{N}\|_{2,1} \text{ s.t. } \mathbf{D} = \mathbf{DR} + \mathbf{N}, \mathbf{R} = \mathbf{X}. \quad (22)$$

Then switch the Equation (22) to an unconstraint problem, which is the Lagrange function. The formula is as follows:

$$\begin{aligned} L_\beta(\mathbf{X}, \mathbf{R}, \mathbf{N}) = & \|\mathbf{X}\|_* + \omega \|\mathbf{N}\|_{2,1} \\ & + \text{tr}(Y_1^T (\mathbf{D} - \mathbf{DR} - \mathbf{N})) + \text{tr}(Y_2^T (\mathbf{R} - \mathbf{X})) \\ & + \frac{\beta}{2} (\|\mathbf{D} - \mathbf{DR} - \mathbf{N}\|_F^2 + \|\mathbf{R} - \mathbf{X}\|_F^2), \end{aligned} \quad (23)$$

where $\beta > 0$ is the penalty parameter, and β is updated by $\beta = \min(\rho\beta, \max_\beta)$. Y_1 and Y_2 are the Lagrange multipliers.

The ADM method is used to solve the Equation (23) [39]. The ADM is a simple method to solve the decomposable convex optimization problem, especially in solving large-scale problems. The update iterations for ADM are as follows:

$$\begin{cases} \mathbf{X}^{k+1} = \arg \min_{\mathbf{X}} L(\mathbf{X}, \mathbf{R}^k, \mathbf{N}^k, \beta), \\ \mathbf{R}^{k+1} = \arg \min_{\mathbf{R}} L(\mathbf{X}^{k+1}, \mathbf{R}, \mathbf{N}^k, \beta), \\ \mathbf{N}^{k+1} = \arg \min_{\mathbf{N}} L(\mathbf{X}^{k+1}, \mathbf{R}^{k+1}, \mathbf{N}, \beta). \end{cases} \quad (24)$$

Based on the singular value shrinkage operator [40], \mathbf{X}^{k+1} and \mathbf{N}^{k+1} are represented

as follows:

$$\mathbf{X}^{k+1} = D_{\frac{1}{\beta}}(\mathbf{R} + \frac{Y_2}{\beta}) = \operatorname{argmin} \frac{1}{\beta} \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - (\mathbf{R} + \frac{Y_2}{\beta})\|_F^2, \quad (25)$$

$$\mathbf{N}^{k+1} = D_{\frac{\omega}{\beta}}(\mathbf{D} - \mathbf{D}\mathbf{R} + \frac{Y_1}{\beta}) = \operatorname{argmin} \frac{\omega}{\beta} \|\mathbf{N}\|_{2,1} + \frac{1}{2} \|\mathbf{N} - (\mathbf{D} - \mathbf{D}\mathbf{R} + \frac{Y_1}{\beta})\|_F^2, \quad (26)$$

yet the minimization of \mathbf{R} is a least squares problem, and its normal equation is as

follows:

$$\mathbf{R} = (\mathbf{I} + \mathbf{D}^T \mathbf{D})^{-1} (\mathbf{D}^T \mathbf{D} - \mathbf{D}^T \mathbf{N} + \mathbf{X} + \frac{\mathbf{D}^T Y_1 - Y_2}{\beta}), \quad (27)$$

where $\mathbf{I} = \mathbf{D}\mathbf{D}^T$ is widely used in matrix completion.

Then \mathbf{X} , \mathbf{R} and \mathbf{N} are updated by changing the Lagrange multipliers Y_1 and Y_2 .

Moreover, Y_1 and Y_2 can be obtained by the following formulas:

$$Y_1 = Y_1 + \beta(\mathbf{D} - \mathbf{D}\mathbf{R} - \mathbf{E}), \quad (28)$$

$$Y_2 = Y_2 + \beta(\mathbf{R} - \mathbf{X}). \quad (29)$$

Finally, we can get the final low-rank matrix \mathbf{R}^* and sparse matrix \mathbf{N}^* until the convergence conditions $\|\mathbf{D} - \mathbf{D}\mathbf{R} - \mathbf{N}\|_{\infty} < \varepsilon$ and $\|\mathbf{R} - \mathbf{X}\|_{\infty} < \varepsilon$ are satisfied. Here, ε is an extremely low number (set as 1×10^{-8} in this paper). As mentioned above, the refined matrix \mathbf{R}^* and noise matrix \mathbf{N}^* can be used to describe a completed matrix in the form of $\mathbf{D} \times \mathbf{R}^*$ or $\mathbf{D} - \mathbf{N}^*$. The specific process of matrix completion is shown in Fig.5.

Based on the above matrix completion method, the disease semantic similarity matrix $\mathbf{D}\mathbf{S}$ and miRNA functional similarity matrix $\mathbf{M}\mathbf{F}$ are used as input matrices to replace matrix \mathbf{D} , so that we can obtain two refined similarity matrices $\mathbf{C}\mathbf{D}$ and $\mathbf{C}\mathbf{M}$, respectively.

The algorithm of Matrix completion is summarized in Algorithm 1 and listed as follows:

Algorithm 1: Matrix completion

Input: an incomplete data matrix \mathbf{D} and $\omega \in (0,1)$

Output: complete matrix $\mathbf{D} \times \mathbf{R}^*$

Initialization: $D = 0, N = 0, Y_1 = 0, Y_2 = 0, \beta = 10^{-4}, \max_{\beta} = 10^{10}, \rho = 1.1, \varepsilon = 10^{-8}$

Repeat:

1. Update \mathbf{X} , \mathbf{R} and \mathbf{N} using Eq.(25), Eq.(26) and Eq.(27), respectively
2. Update the multiplier Y_1 and Y_2 using Eq.(28) and Eq.(29), respectively
3. Update parameter β by: $\beta = \min(\rho\beta, \max_{\beta})$
4. Compare the convergence conditions: $\|\mathbf{D} - \mathbf{D}\mathbf{R} - \mathbf{N}\|_{\infty} < \varepsilon$, $\|\mathbf{R} - \mathbf{X}\|_{\infty} < \varepsilon$

Until convergence

Return: $\mathbf{D} \times \mathbf{R}^*$ or $\mathbf{D} - \mathbf{N}^*$

Similarity Information Integrations

Subsequent work is to integrate the completed matrix with existing similarity matrices.

Since similarity information integrations of diseases and miRNAs are similar, Fig.6 only shows the process for integration of miRNA similarity.

The specific integration formulas are as follows:

$$\mathbf{IMS}(i, j) = \begin{cases} \frac{\mathbf{CM}(i, j) + \mathbf{GM}(i, j)}{2}, & \text{if } \mathbf{MF}(i, j) = 0, \\ \frac{\mathbf{GM}(i, j) + \mathbf{CM}(i, j) + \mathbf{MF}(i, j)}{3}, & \text{otherwise,} \end{cases} \quad (30)$$

$$\mathbf{IDS}(i, j) = \begin{cases} \frac{\mathbf{CD}(i, j) + \mathbf{GD}(i, j)}{2}, & \text{if } \mathbf{DS}(i, j) = 0, \\ \frac{\mathbf{GD}(i, j) + \mathbf{CD}(i, j) + \mathbf{DS}(i, j)}{3}, & \text{otherwise,} \end{cases} \quad (31)$$

WKNKN

WKNKN can be thought of as a voting or integration method: some potential classifiers (nearest neighbors) are aggregated by a (weight) majority vote, the results of which are used for prediction [41].

In this paper, \mathbf{MD} expresses the miRNA-disease association matrix, which only represents the association between the miRNA and the disease verified by human experiment at the current stage. And we simply stipulate that if the miRNA is associated with the disease, $\mathbf{MD}(m_i, d_j)$ will be set to 1. However, there are still many unknown miRNAs and diseases in the world, and whether they can be used as a bridge between existing miRNAs and diseases or not are still unknown. Maybe existing miRNAs are correlated with existing diseases through these unknown miRNAs, so the \mathbf{MD} regulation is obviously inappropriate.

Therefore, by estimating these unknown conditions through the correlation of its known neighbors, the WKNKN method preprocesses the matrix \mathbf{MD} to get the pre-processed matrix of \mathbf{MD} (\mathbf{PMD}). If $\mathbf{MD}(m_i, d_j) = 0$, WKNKN will give $\mathbf{MD}(m_i, d_j)$ a value from 0 to 1 according to the corresponding similar information of miRNAs and diseases. The specific process of WKNKN is shown in Fig.7.

MCCMF for MiRNA-Disease Association Prediction

The CMF method proposed by Shen *et al.* that can effectively predict the potential interactions between miRNAs and diseases [45]. In this study, the idea of the CMF method is used to predict the miRNA-disease association. The specific steps of CMF are as follows: firstly, the input miRNA-disease association matrix \mathbf{PMD} is decomposed into two low-rank matrices \mathbf{A} and \mathbf{B} by using the singular value

decomposition.

$$\begin{aligned}
[\mathbf{U}, \mathbf{S}, \mathbf{V}] &= SVD(\mathbf{PMD}, k), \\
\mathbf{A} &= \mathbf{U} \mathbf{S}_k^{\frac{1}{2}}, \\
\mathbf{B} &= \mathbf{V} \mathbf{S}_k^{\frac{1}{2}},
\end{aligned} \tag{32}$$

where \mathbf{U} and \mathbf{V} is the unitary matrix. \mathbf{S} is a negative real diagonal matrix, and there are k singular values on the diagonal.

Secondly, we write the objection function of MCCMF according to the idea of CMF, as follows:

$$\begin{aligned}
\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{PMD} - \mathbf{A}\mathbf{B}^T\|_F^2 &+ \lambda_l (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2) \\
&+ \lambda_m \|\mathbf{IMS} - \mathbf{A}\mathbf{A}^T\|_F^2 \\
&+ \lambda_d \|\mathbf{IDS} - \mathbf{B}\mathbf{B}^T\|_F^2.
\end{aligned} \tag{33}$$

Here, $\|\cdot\|_F$ is the Frobenius norm to ensure that the feature vectors of similar miRNAs and similar diseases are similar. λ_l , λ_m and λ_d are positive parameters, which are determined by the 5-fold cross validation, and $\lambda_l \in \{2^{-2}, 2^{-1}, 2^0, 2^1\}$, $\lambda_m / \lambda_d \in \{2^{-3}, 2^{-2}, 2^{-1}, 2^0, 2^1, 2^2, 2^3, 2^4, 2^5\}$.

Thirdly, we use L to represent the equation (33), and derive two alternative update rules by setting $\partial L / \partial \mathbf{A} = 0$ and $\partial L / \partial \mathbf{B} = 0$.

$$\mathbf{A} = (\mathbf{PMD} * \mathbf{B} + \lambda_m \mathbf{IMS} * \mathbf{A})(\mathbf{B}^T \mathbf{B} + \lambda_l \mathbf{I}_k + \lambda_m \mathbf{A}^T \mathbf{A})^{-1}, \tag{34}$$

$$\mathbf{B} = (\mathbf{PMD}^T * \mathbf{A} + \lambda_d \mathbf{IDS} * \mathbf{B})(\mathbf{A}^T \mathbf{A} + \lambda_l \mathbf{I}_k + \lambda_d \mathbf{B}^T \mathbf{B})^{-1}, \tag{35}$$

where \mathbf{I}_k is the $k \times k$ identity matrix.

Finally, we update \mathbf{A} and \mathbf{B} iteratively until they converge to get the final \mathbf{A} and \mathbf{B} . By $\mathbf{A} * \mathbf{B}^T$, the prediction matrix for miRNA-disease associations is obtained. The detail process of MCCMF can be seen in Fig.8.

The algorithm of CMF is summarized in Algorithm 2 and listed as follows:

Algorithm 2: CMF:

Input: pre-processed matrix **PMD**, integration of miRNA similarity **IMS** and integration of disease similarity **IDS**

Output: prediction score matrix $\mathbf{A} * \mathbf{B}^T$

Initialization: $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{SVD}(\mathbf{Y}, k)$, $\mathbf{A} = \mathbf{US}_k^{1/2}$, $\mathbf{B} = \mathbf{VS}_k^{1/2}$

Repeat:

Update **A** and **B** using Eq.(34) and Eq.(35), respectively

Until convergence

Return: $\mathbf{A} * \mathbf{B}^T$

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Availability of data and materials

The datasets that support the findings of this study are available in

<https://github.com/cuizhensdws>.

Competing interests

The authors declare that they have no competing interests

Funding

Publication costs are funded by the National Science Foundation of China under grant

Nos. 61872220, and 61702299.

Authors' contributions

TRW and MMY jointly contributed to the design of the study. TRW designed and implemented the MCCMF method, performed the experiments, and drafted the manuscript. XZK participated in the design of the study and performed the statistical analysis. YLG contributed to the data analysis. CNJ and JXL gave computational advice for the project and participated in designing evaluation criteria. All authors read and approved the final manuscript.

Acknowledgements

Not applicable

Reference

1. Alshalalfa M, Alhadj R: **Using context-specific effect of miRNAs to identify functional associations between miRNAs and gene signatures.** *BMC Bioinformatics* 2013, **14**(12):S1.
2. Lee RC, Feinbaum RL, Ambros V: **The C. elegans heterochronic gene lin-4 encodes small RNAs with antisense complementarity to lin-14.** *Cell* 1993, **75**(5):843-854.
3. Kozomara A, Griffiths-Jones S: **miRBase: annotating high confidence microRNAs using deep sequencing data.** *Nucleic Acids Research* 2013, **42**(D1):D68-D73.
4. Cheng AM, Byrom MW, Shelton J, Ford LP: **Antisense inhibition of human miRNAs and indications for an involvement of miRNA in cell growth and apoptosis.** *Nucleic Acids Research* 2005, **33**(4):1290-1297.
5. Miska EA: **How microRNAs control cell division, differentiation and death.** *Current Opinion in Genetics & Development* 2005, **15**(5):563-568.
6. Bartel DP: **MicroRNAs: Target Recognition and Regulatory Functions.** *Cell* 2009, **136**(2):215-233.
7. Xu P, Guo M, Hay BA: **MicroRNAs and the regulation of cell death.** *Trends in Genetics* 2004, **20**(12):617-624.
8. Lewis BP, Burge CB, Bartel DP: **Conserved Seed Pairing, Often Flanked by Adenosines, Indicates that Thousands of Human Genes are MicroRNA Targets.** *Cell* 2005, **120**(1):15-20.
9. Lu M, Zhang Q, Deng M, Miao J, Guo Y, Gao W, Cui Q: **An Analysis of Human MicroRNA and Disease Associations.** *PLOS ONE* 2008, **3**(10):e3420.
10. Calin GA, Dumitru CD, Shimizu M, Bichi R, Zupo S, Noch E, Aldler H,

-
- Rattan S, Keating M, Rai K *et al*: **Frequent deletions and down-regulation of micro- RNA genes miR15 and miR16 at 13q14 in chronic lymphocytic leukemia.** *Proceedings of the National Academy of Sciences of the United States of America* 2002, **99**(24):15524-15529.
11. Wu C, Li M, Hu C, Duan H: **Clinical significance of serum miR-223, miR-25 and miR-375 in patients with esophageal squamous cell carcinoma.** *Molecular Biology Reports* 2014, **41**(3):1257-1266.
 12. Zhang X, Zhang X, Wang T, Wang L, Zhijun T, Wei W, Yan B, Zhao J, Wu K, Yang A-G *et al*: **MicroRNA-26a is a key regulon that inhibits progression and metastasis of c-Myc/EZH2 double high advanced hepatocellular carcinoma.** *Cancer Letters* 2018, **426**.
 13. Wu Z-s, Wu Q, Wang C-q, Wang X-n, Huang J, Zhao J-j, Mao S-s, Zhang G-h, Xu X-c, Zhang N: **miR-340 Inhibition of Breast Cancer Cell Migration and Invasion Through Targeting of Oncoprotein c-Met.** *Cancer* 2011, **117**(13):2842-2852.
 14. Yang Z, Ren F, Liu C, he S, Sun G, Gao Q, Yao L, Zhang Y, Miao R, Cao Y *et al*: **DbDEMC: A database of differentially expressed miRNAs in human cancers.** *BMC genomics* 2010, **11** Suppl 4:S5.
 15. Huang Z, Shi J, Gao Y, Cui C, Zhang S, Li J, Zhou Y, Cui Q: **HMDD v3.0: a database for experimentally supported human microRNA–disease associations.** *Nucleic Acids Research* 2019, **47**(D1):D1013-D1017.
 16. Jiang Q, Wang Y, Hao Y, Juan L, Teng M, Zhang X, Li M, Wang G, Liu Y: **miR2Disease: a manually curated database for microRNA deregulation in human disease.** *Nucleic acids research* 2008, **37**:D98-104.
 17. Chen X, Xie D, Zhao Q, You Z-H: **MicroRNAs and complex diseases: from experimental results to computational models.** *Briefings in Bioinformatics* 2019, **20**(2):515-539.
 18. Zou Q, Li J, Song L, Zeng X, Wang G: **Similarity computation strategies in the microRNA-disease network: A survey.** *Briefings in functional genomics* 2015, **15**(18):55-64.
 19. Jiang Q, Hao Y, Wang G, Juan L, Zhang T, Teng M, Liu Y, Wang Y: **Prioritization of disease microRNAs through a human phenome-microRNAome network.** *BMC Systems Biology* 2010, **4**(1):S2.
 20. Li X, Wang Q, Zheng Y, Lv S, Ning S, Sun J, Huang T, Zheng Q, Ren H, Xu J *et al*: **Prioritizing human cancer microRNAs based on genes' functional consistency between microRNA and cancer.** *Nucleic Acids Research* 2011, **39**(22):e153-e153.
 21. Xiao Q, Luo J, Liang C, Cai J, Ding P: **A graph regularized non-negative matrix factorization method for identifying microRNA-disease associations.** *Bioinformatics* 2017, **34**(2):239-248.
 22. Chen X, Yin J, Qu J, Huang L: **MDHGI: Matrix Decomposition and Heterogeneous Graph Inference for miRNA-disease association prediction.** *PLOS Computational Biology* 2018, **14**(8):e1006418.
 23. Mørk S, Pletscher-Frankild S, Palleja A, Gorodkin J, Jensen L:

-
- Protein-driven inference of miRNA-disease associations.** *Bioinformatics (Oxford, England)* 2013, **30**(3).
24. Chen H, Zhang Z: **Similarity-based methods for potential human microRNA-disease association prediction.** *BMC medical genomics* 2013, **6**:12.
 25. Gao M-M, Cui Z, Gao Y-L, Liu J-X, Zheng C-H: **Dual-network sparse graph regularized matrix factorization for predicting miRNA–disease associations.** *Molecular Omics* 2019, **15**(2):130-137.
 26. Gao Y-L, Cui Z, Liu J-X, Wang J, Zheng C-H: **NPCMF: Nearest Profile-based Collaborative Matrix Factorization method for predicting miRNA-disease associations.** *BMC Bioinformatics* 2019, **20**(1):353.
 27. Yin M-M, Cui Z, Gao M-M, Liu J-X, Gao Y-L: **LWPCMF: Logistic Weighted Profile-based Collaborative Matrix Factorization for Predicting MiRNA-Disease Associations.** *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 2019, **PP**:1-1.
 28. Chen H, Zhang Z, Feng D: **Prediction and interpretation of miRNA-disease associations based on miRNA target genes using canonical correlation analysis.** *BMC Bioinformatics* 2019, **20**(1):404.
 29. Xu J, Li C-X, Lv J-Y, Li Y-S, Xiao Y, Shao T-T, Huo X, Li X, Zou Y, Han Q-L *et al*: **Prioritizing Candidate Disease miRNAs by Topological Features in the miRNA Target–Dysregulated Network: Case Study of Prostate Cancer.** *Molecular Cancer Therapeutics* 2011, **10**(10):1857.
 30. Chen H, Zhang Z: **Prediction of Associations between OMIM Diseases and MicroRNAs by Random Walk on OMIM Disease Similarity Network.** *TheScientificWorldJournal* 2013, **2013**:204658.
 31. Chen X, Yan CC, Zhang X, Li Z, Deng L, Zhang Y, Dai Q: **RBMMMDA: predicting multiple types of disease-microRNA associations.** *Scientific Reports* 2015, **5**:13877.
 32. Chen X, Huang L: **LRSSLMDA: Laplacian Regularized Sparse Subspace Learning for MiRNA-Disease Association prediction.** *PLOS Computational Biology* 2017, **13**(12):e1005912.
 33. Chen X, Wang L, Qu J, Guan N-N, Li J-Q: **Predicting miRNA–disease association based on inductive matrix completion.** *Bioinformatics* 2018, **34**(24):4256-4265.
 34. Chen X, Xie D, Wang L, Zhao Q, You Z-H, Liu H: **BNPMDA: Bipartite Network Projection for MiRNA–Disease Association prediction.** *Bioinformatics* 2018, **34**(18):3178-3186.
 35. Chen X, Zhu C-C, Yin J: **Ensemble of decision tree reveals potential miRNA-disease associations.** *PLoS Computational Biology* 2019, **15**(7):e1007209.
 36. Ding X, Xia J-F, Wang Y-T, Wang J, Zheng C-H: **Improved Inductive Matrix Completion Method for Predicting MicroRNA-Disease Associations.** In.; 2019: 247-255.
 37. Li J, Zhang S, Liu T, Ning C, Zhang Z, Zhou W: **Neural inductive matrix**

-
- completion with graph convolutional networks for miRNA-disease association prediction. *Bioinformatics* 2020, **36**(8):2538-2546.
38. Lin Z, Chen M, Ma Y: **The Augmented Lagrange Multiplier Method for Exact Recovery of Corrupted Low-Rank Matrices.** *Mathematical Programming* 2010, **9**.
39. Yang J, Yuan X: **Linearized augmented Lagrangian and alternating direction methods for nuclear norm minimization.** *Mathematics of Computation* 2011, **82**.
40. Cai J-F, Candès EJ, Shen Z: **A Singular Value Thresholding Algorithm for Matrix Completion.** *SIAM Journal on Optimization* 2010, **20**:1956-1982.
41. Ezzat A, Zhao P, Wu M, Li X, Kwoh C: **Drug-Target Interaction Prediction with Graph Regularized Matrix Factorization.** *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 2017, **14**(3):646-656.
42. Ezzat A, Wu M, Li X-L, Kwoh C-K: **Computational prediction of drug-target interactions using chemogenomic approaches: an empirical survey.** *Briefings in Bioinformatics* 2018, **20**(4):1337-1357.
43. Xie G, Fan Z, Sun Y, Wu C, Ma L: **WBNPMD: weighted bipartite network projection for microRNA-disease association prediction.** *Journal of Translational Medicine* 2019, **17**:322.
44. Chen X, Yan G-Y: **Semi-supervised learning for potential human microRNA-disease associations inference.** *Scientific Reports* 2014, **4**:5501.
45. Shen Z, Zhang Y-H, Han K, Nandi A, Honig B, Huang D-S: **miRNA-Disease Association Prediction with Collaborative Matrix Factorization.** *Complexity* 2017, **2017**:1-9.
46. Huang Z, Shi J, Gao Y, Cui C, Zhang S, Li J, Zhou Y, Cui Q: **HMDD v3.0: a database for experimentally supported human microRNA-disease associations.** *Nucleic acids research* 2018, **47**.
47. Xie B, Ding Q, Han H, Wu D: **MiRCancer: A microRNA-cancer association database constructed by text mining on literature.** *Bioinformatics (Oxford, England)* 2013, **29**.
48. Tazawa H, Kagawa S, Fujiwara T: **MicroRNAs as potential target gene in cancer gene therapy of gastrointestinal tumors.** *Expert opinion on biological therapy* 2011, **11**:145-155.
49. Montoya V, Fan H, Bryar P, Weinstein J, Mets M, Feng G, Martin J, Martin A, Jiang H, Laurie N: **Novel miRNA-31 and miRNA-200a-Mediated Regulation of Retinoblastoma Proliferation.** *PloS one* 2015, **10**:e0138366.
50. Zhang X, Liu S, Hu T, Liu S, He Y, Sun S: **Up-Regulated MicroRNA-143 Transcribed by Nuclear Factor kappa B Enhances Hepatocarcinoma Metastasis by Repressing Fibronectin Expression.** *Hepatology (Baltimore, Md)* 2009, **50**:490-499.
51. Li Y, Qiu C, Tu J, Geng B, Yang J, Jiang T, Cui Q: **HMDD v2.0: A database for experimentally supported human microRNA and disease associations.** *Nucleic acids research* 2013, **42**.
52. Wang D, Wang J, Lu M, Song F, Cui Q: **Inferring the human microRNA**

-
- functional similarity and functional network based on microRNA-associated diseases.** *Bioinformatics (Oxford, England)* 2010, **26**:1644-1650.
53. Chen H, Guo R, Li G, Zhang W, Zhang Z: **Comparative analysis of similarity measurements in miRNAs with applications to miRNA-disease association predictions.** *BMC Bioinformatics* 2020, **21**(1):176.
54. Xuan P, Han K, Guo M, Guo Y, Li J, Ding J, Liu Y, Dai Q, Li J, Teng Z *et al*: **Prediction of microRNAs Associated with Human Diseases Based on Weighted k Most Similar Neighbors.** *PloS one* 2013, **8**:e70204.
55. van Laarhoven T, Nabuurs SB, Marchiori E: **Gaussian interaction profile kernels for predicting drug–target interaction.** *Bioinformatics* 2011, **27**(21):3036-3043.
56. Chen X, Yan G-Y: **Novel human lncRNA–disease association inference based on lncRNA expression profiles.** *Bioinformatics* 2013, **29**(20):2617-2624.
57. Sheng-Peng Y, Liang C, Xiao Q, Li GH, Ding P, Luo JW: **MCLPMDA: A novel method for miRNA-disease association prediction based on matrix completion and label propagation.** *Journal of Cellular and Molecular Medicine* 2018, **23**.

Tables

Table 1 AUC results of cross validation experiments

Methods	Gold Standard Dataset
RLSMDA	0.8389(0.0006)
GRNMF	0.869(0.00023)
CMF	0.8697(0.0011)
WBNPMD	0.9173(0.0005)
MCCMF	0.9569(0.0005)

Table 2 Predicted miRNAs for Gastrointestinal Neoplasms

Rank	MiRNA	Evidence	Rank	MiRNA	Evidence
1	has-mir-1	known	21	has-let-7a	known
2	has-mir-22	known	22	has-mir-152	known
3	has-mir-200	known	23	has-mir-497	known
4	has-mir-9	known	24	has-mir-21	HMDD v3.0
5	has-mir-221	known	25	has-mir-375	known
6	has-mir-146a	known	26	has-mir-107	known
7	has-mir-133b	known	27	has-mir-18b	known
8	has-mir-200c	known	28	has-mir-494	known
9	has-mir-200a	known	29	has-mir-150	miRCancer
10	has-mir-7	known	30	has-mir-208a	known
11	has-mir-200b	known	31	has-mir-98	known
12	has-mir-222	known	32	has-mir-141	miRCancer
13	has-mir-126	known	33	has-let-7g	known
14	has-mir-196a	known	34	has-mir-184	unconfirmed
15	has-mir-142	known	35	has-mir-210	miRCancer
16	has-mir-124	known	36	has-mir-486	HMDD v3.0
17	has-mir-148a	known	37	has-mir-338	known
18	has-mir-451a	known	38	has-mir-27a	miRCancer
19	has-mir-31	known	39	has-mir-146b	HMDD v3.0
20	has-mir-451	known	40	has-let-7c	unconfirmed

Table 3 Predicted microbes for Retinoblastoma

Rank	MiRNA	Evidence	Rank	MiRNA	Evidence
1	has-mir-1	known	31	has-mir-32	unconfirmed
2	has-mir-9	known	32	has-mir-200	unconfirmed
3	has-mir-17	known	33	has-mir-192	known
4	has-mir-20a	known	34	has-mir-513b	known
5	has-mir-18a	known	35	has-mir-135b	known
6	has-mir-29c	known	36	has-mir-513c	known
7	has-mir-92a	known	37	has-mir-22	miRCancer
8	has-let-7d	known	38	has-mir-31	HMDD v3.0

9	has-let-7f	known	39	has-mir-513a	known
10	has-let-7g	known	40	has-mir-30c	known
11	has-let-7a	known	41	has-mir-491	known
12	has-mir-19b	known	42	has-mir-135a	unconfirmed
13	has-let-7b	known	43	has-mir-125b	HMDD v3.0
14	has-mir-29b	known	44	has-mir-7	unconfirmed
15	has-let-7e	known	45	has-mir-181a	unconfirmed
16	has-let-7i	known	46	has-mir-223	unconfirmed
17	has-mir-124	known	47	has-mir-210	unconfirmed
18	has-let-7c	known	48	has-mir-376a	known
19	has-mir-19a	known	49	has-mir-30a	unconfirmed
20	has-mir-92	known	50	has-mir-145	miRCancer
21	has-mir-181b	known	51	has-mir-34b	HMDD v3.0
22	has-mir-34a	known	52	has-mir-155	unconfirmed
23	has-mir-29a	known	53	has-mir-133a	unconfirmed
24	has-mir-181	known	54	has-mir-137	unconfirmed
25	has-mir-24	known	55	has-mir-146b	unconfirmed
26	has-mir-142	known	56	has-mir-150	unconfirmed
27	has-mir-10b	known	57	has-mir-126	unconfirmed
28	has-mir-34c	known	58	has-mir-18b	unconfirmed
29	has-mir-125a	known	59	has-mir-221	miRCancer
30	has-mir-21	HMDDv3.0 /miRCancer	60	has-mir-373	HMDD v3.0

Table 4 Predicted microbes for Hepatoblastoma

Rank	MiRNA	Evidence	Rank	MiRNA	Evidence
1	has-mir-1	known	11	has-mir-31	unconfirmed
2	has-mir-21	known	12	has-mir-126	dbDEMC
3	has-mir-150	known	13	has-mir-146a	HMDD v3.0
4	has-mir-199a	HMDD v3.0	14	has-mir-125b	unconfirmed
5	has-mir-143	HMDD v3.0	15	has-mir-148a	known
6	has-mir-145	known	16	has-mir-22	dbDEMC
7	has-mir-199b	known	17	has-mir-210	HMDD v3.0
8	has-mir-214	known	18	has-mir-138	unconfirmed
9	has-mir-125a	known	19	has-mir-133a	unconfirmed
10	has-mir-9	unconfirmed	20	has-mir-122	HMDD v3.0

Figures

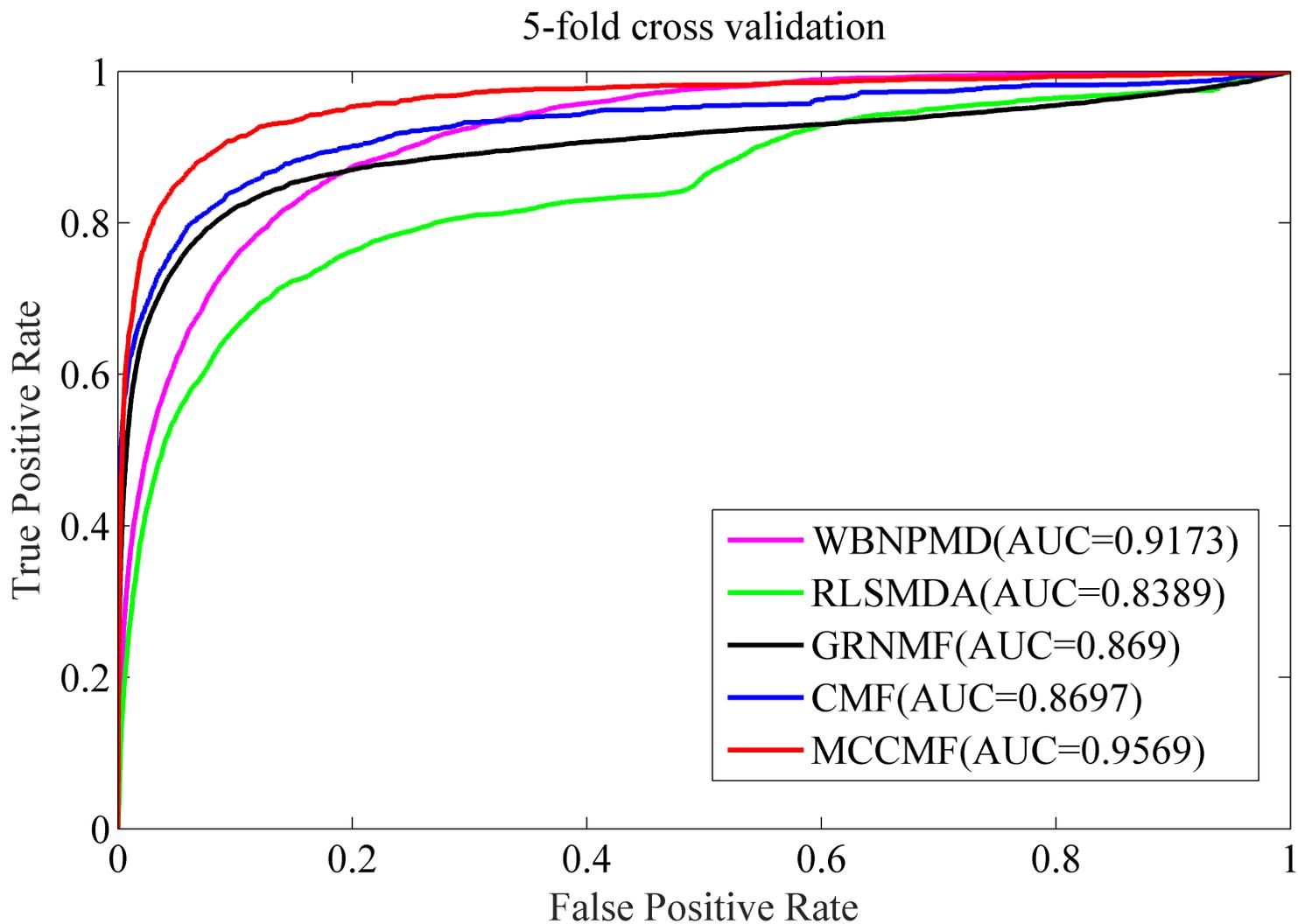


Figure 1

The ROC curves are drawn in Fig.1

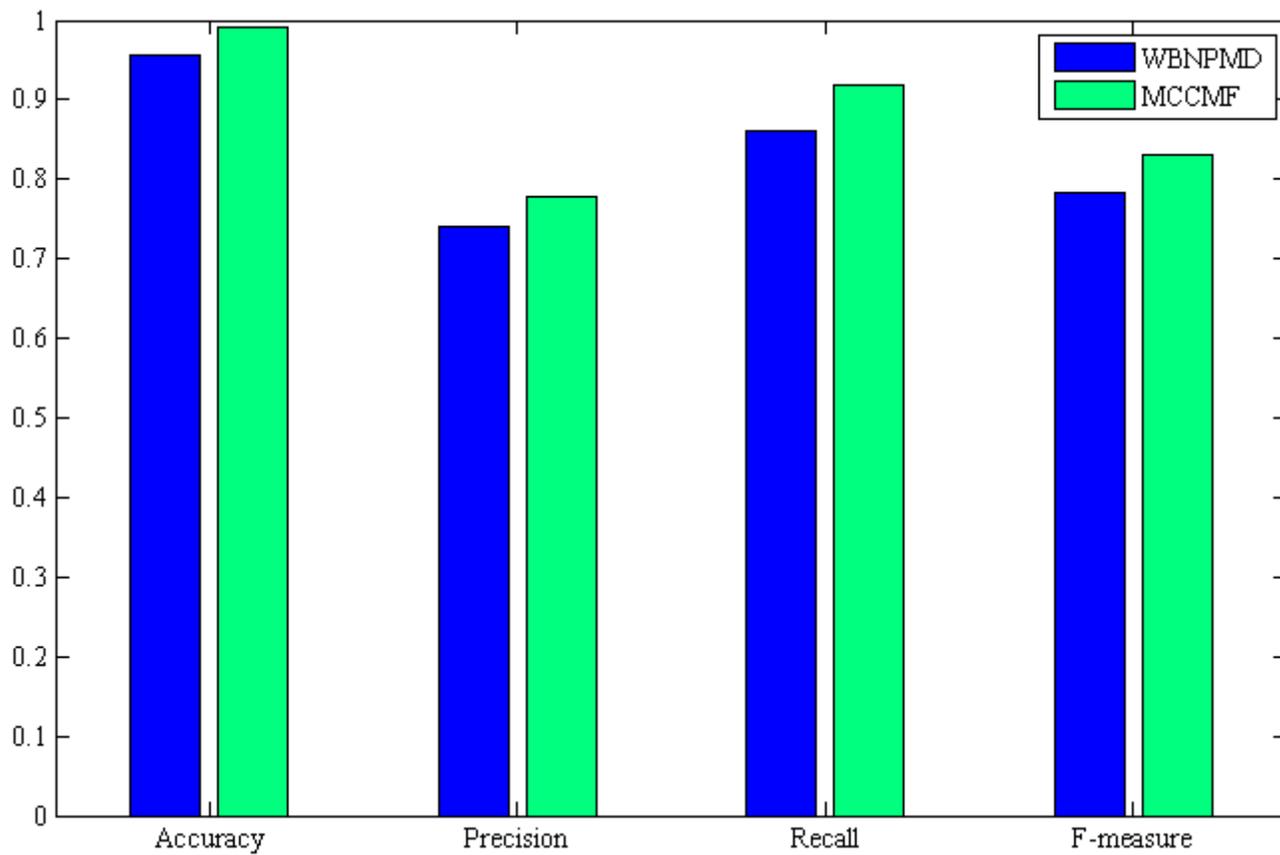


Figure 2

WBNPMD with higher AUC value is selected for comparison with MCCMF, and accuracy, precision, recall and f-measure are presented as a bar graph in Fig.2.

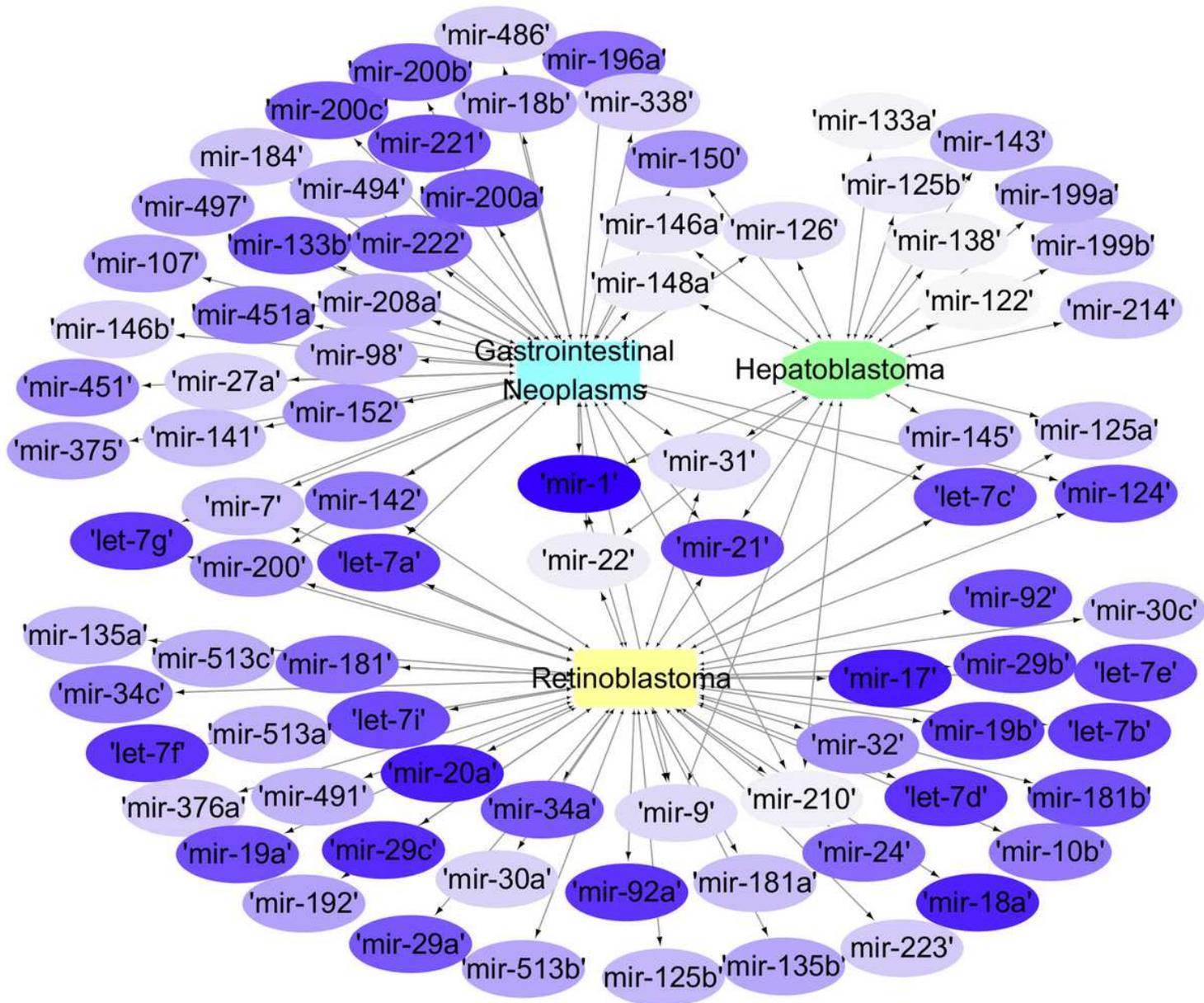


Figure 3

we used Cytoscape software to map the prediction network of these three diseases (Fig.3)

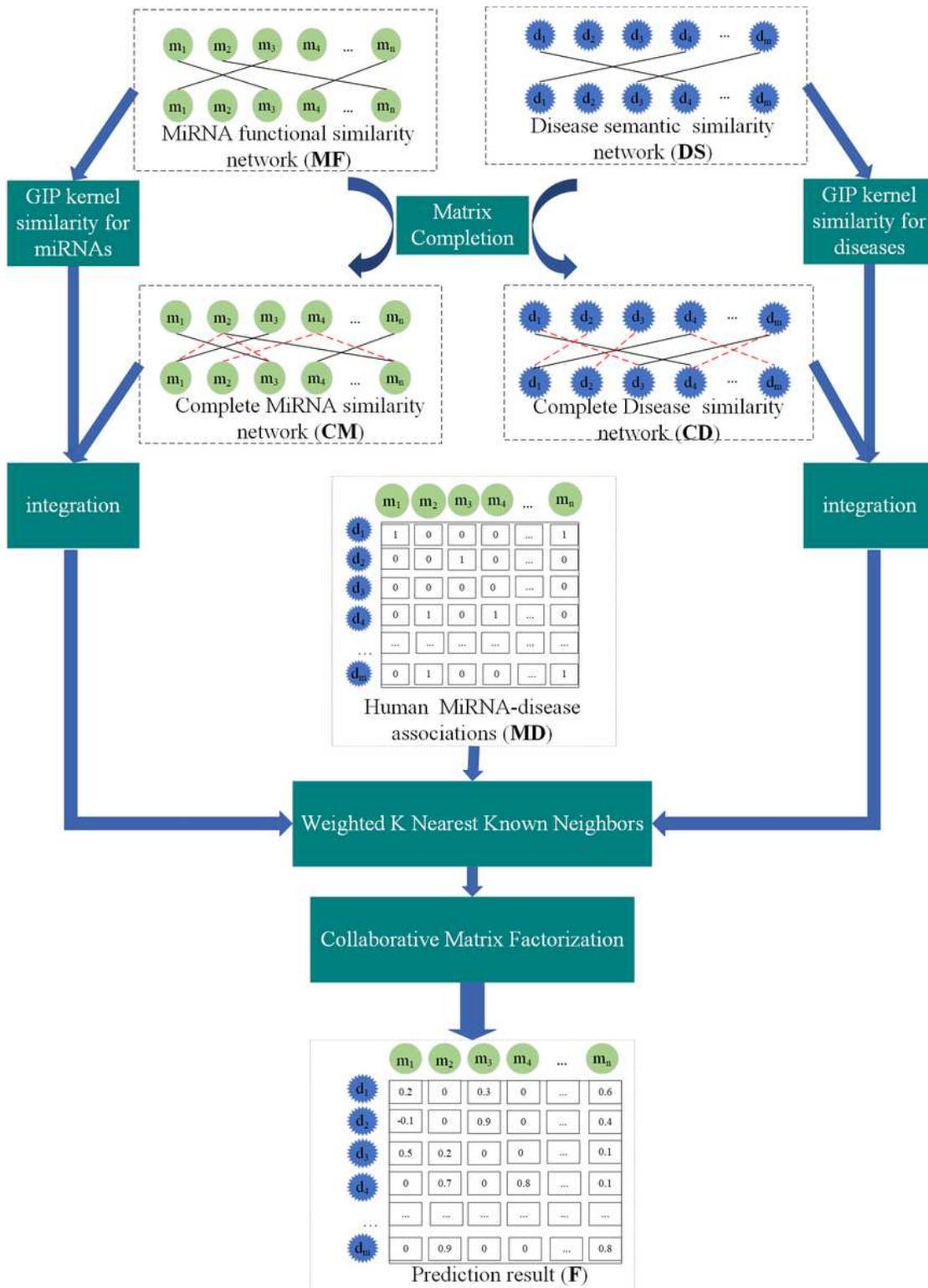


Figure 4

The Collaborative Matrix Factorization is used to predict the association of miRNA-disease. Fig. 4 shows the complete process for MCCMF.

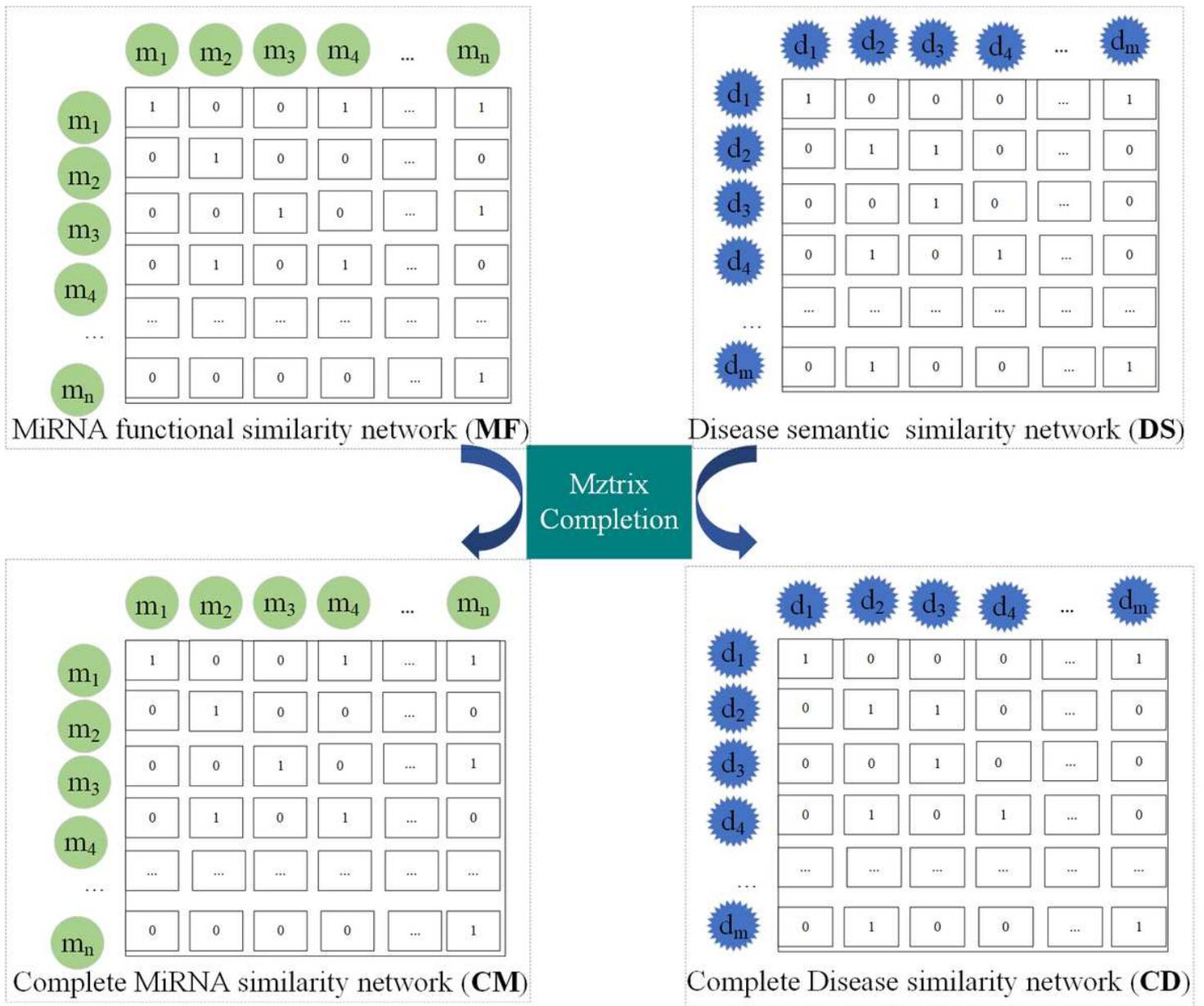


Figure 5

The process of matrix completion

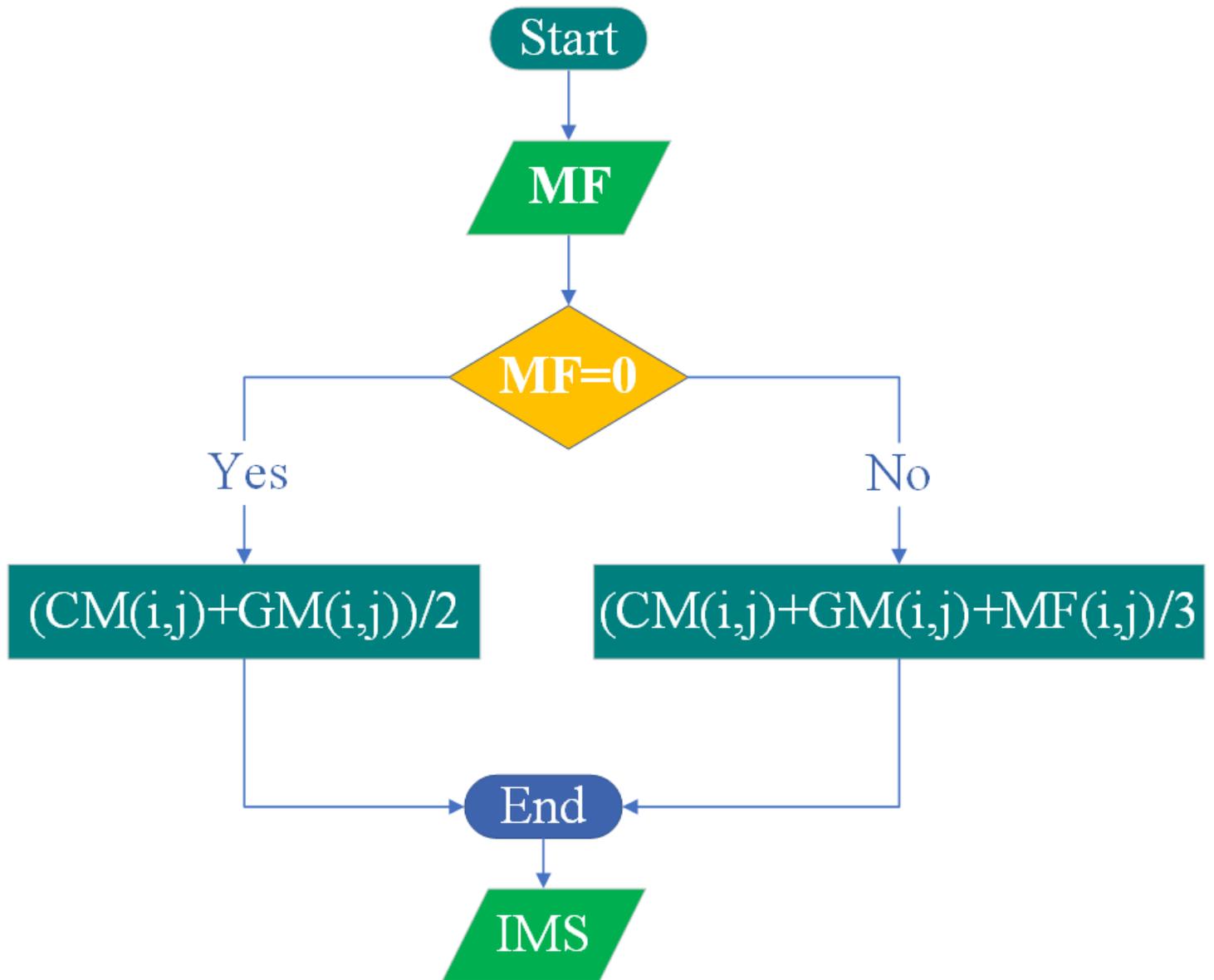


Figure 6

Flowchart of the process for integration of miRNA similarity.

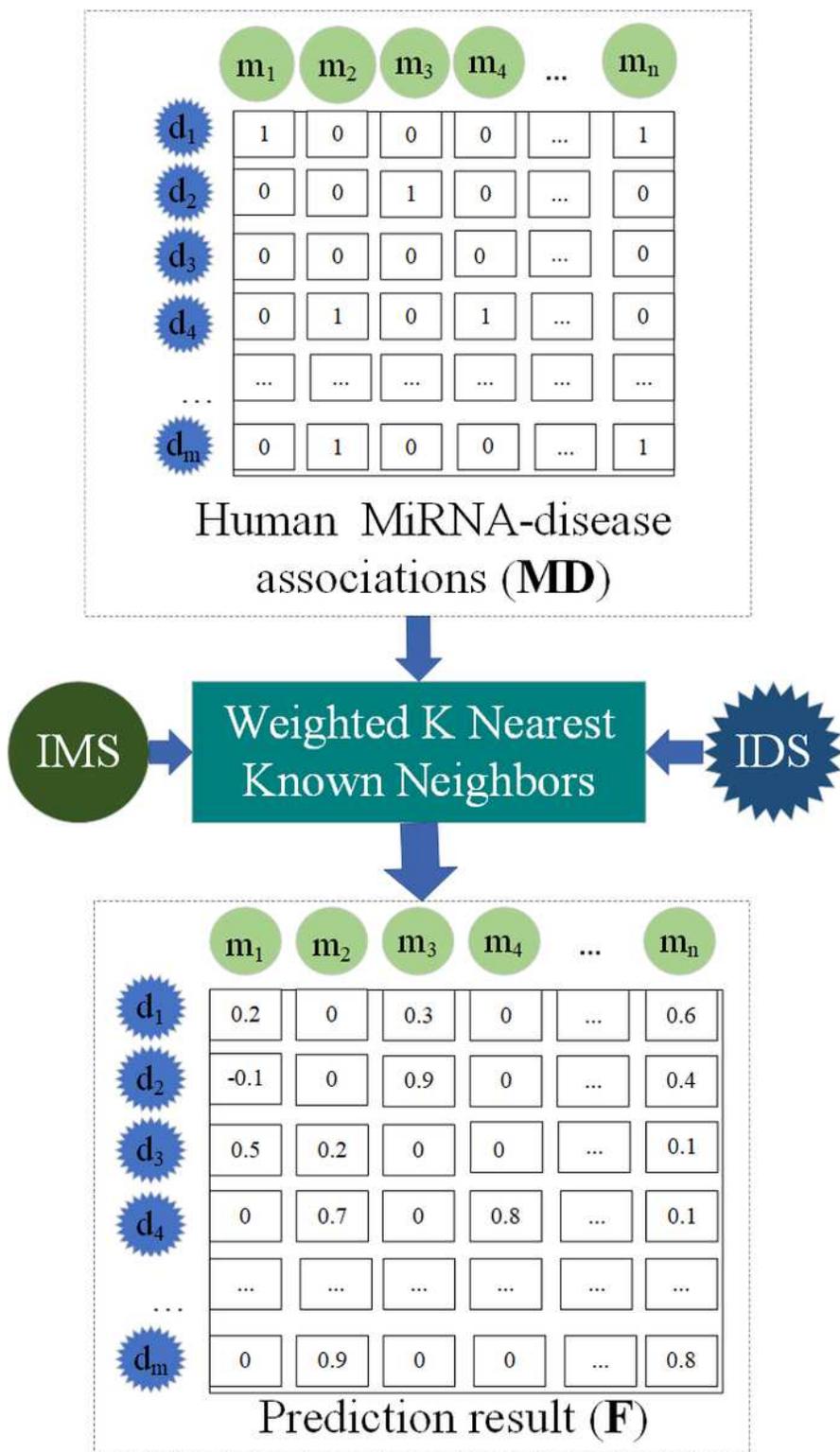


Figure 7

The specific process of WKNKN is shown in Fig.7

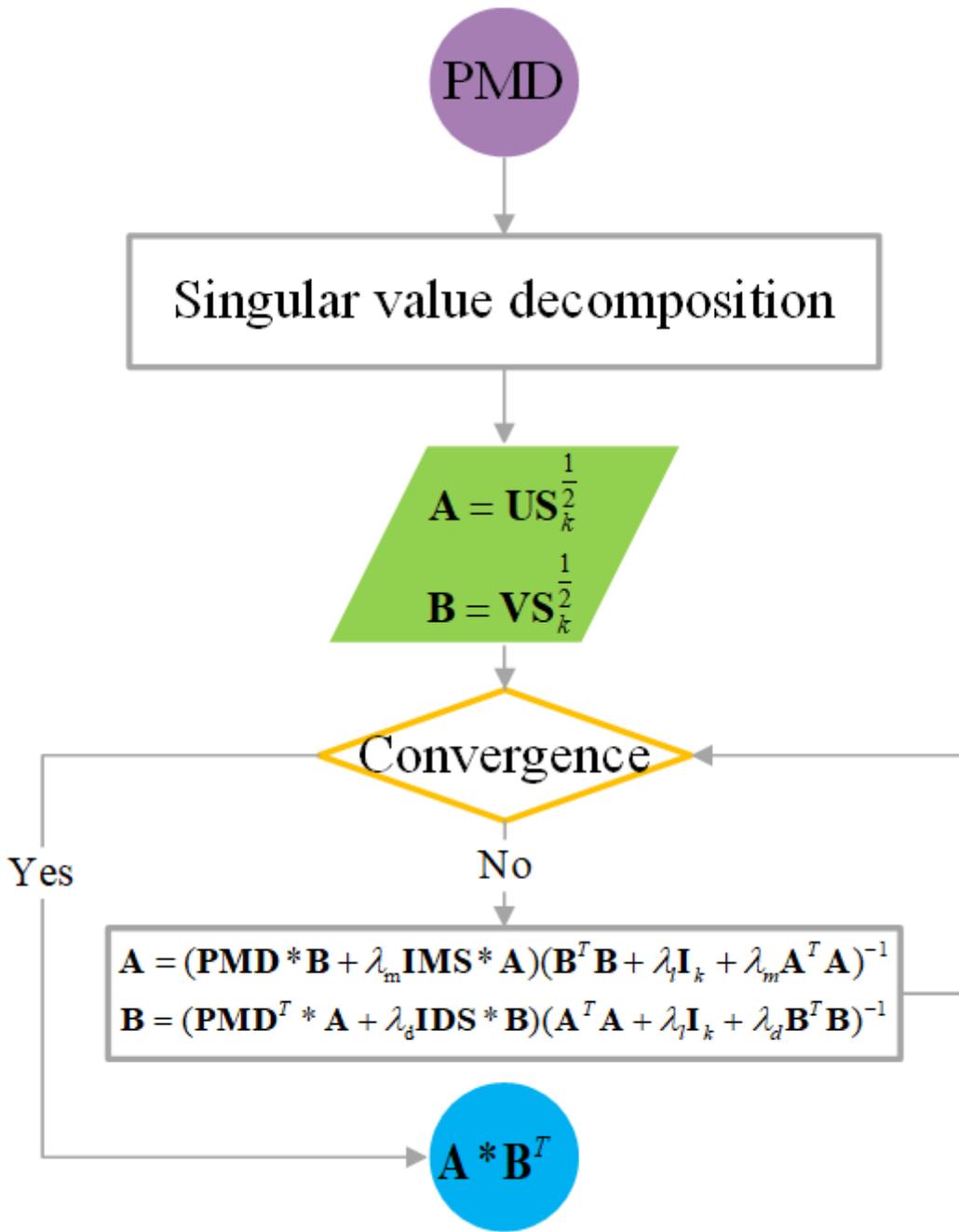


Figure 8

The detail process of MCCMF can be seen in Fig.8.