

Article

A Fusion Link Prediction Method Based on Limit Theorem

Yiteng Wu *, Hongtao Yu *, Ruiyang Huang, Yingle Li and Senjie Lin

National Digital Switching System Engineering and Technological R&D Center, Zhengzhou 450002, China; 18337176095@139.com (R.H.); lyl7225@163.com (Y.L.); 13215997797@163.com (S.L.)

* Correspondence: wuyiteng1992@163.com (Y.W.); 15937101921@139.com (H.Y.); Tel.: +86-150-9348-9229 (Y.W.)

Received: 3 December 2017; Accepted: 22 December 2017; Published: 28 December 2017

Featured Application: The proposed theory can guide the design of combination methods, and the proposed TLF method can fuse multiple similarity indices in link prediction.

Abstract: The theoretical limit of link prediction is a fundamental problem in this field. Taking the network structure as object to research this problem is the mainstream method. This paper proposes a new viewpoint that link prediction methods can be divided into single or combination methods, based on the way they derive the similarity matrix, and investigates whether there a theoretical limit exists for combination methods. We propose and prove necessary and sufficient conditions for the combination method to reach the theoretical limit. The limit theorem reveals the essence of combination method that is to estimate probability density functions of existing links and nonexistent links. Based on limit theorem, a new combination method, theoretical limit fusion (TLF) method, is proposed. Simulations and experiments on real networks demonstrated that TLF method can achieve higher prediction accuracy.

Keywords: link prediction; combination method; theoretical limit; TLF method

1. Introduction

Limit theory is a basic theoretical issue and has attracted wide interest across many fields. On the 100th anniversary of its foundation, *Science* raised 125 unresolved scientific questions, and many of these issues related to limit theory [1]. Link prediction predicts missing links in current networks and new or dissolution links in future networks [2]. With continuous improvement of link prediction methods and, the theoretical limit of link prediction has attracted considerable research interest [3].

Considering structure or attribute features, link prediction methods based on classification have been proposed by computer science community [4,5]. Subsequently, more insightful methods of network structure, such as similarity based methods [6], have become a focus, these methods pay more attention to the physical meaning. At the same time, similarity index fusion methods are springing up [7,8]. Recent years, with the development of deep learning, some deep features extraction methods have been proposed [9,10], the fusion of structure and attribute information has been attached importance again [11–14]. These methods have strong consistency. We divide link prediction method into single and combination methods, based on whether they use multidimension information, and whether they define the relation of multidimension information directly. For example, single methods, such as RA index [15], which defines the relation of common neighbors and degree of nodes directly; and classification based methods, index fusion methods, fusion of structure and attribute information methods belong to link prediction combination methods.

Most combination methods perform better than single methods that will be fused, and are robust to many network types. However, what is the reason for this improved accuracy and robustness, and is there a theoretical limit for combination methods? This paper proposes the mathematic

description of combination methods, and obtains the necessary and sufficient conditions for theoretical limit. The limit theorem also has important practical application value. It reveals the ultimate goal of combination methods that is to estimate probability density functions of existing links and nonexistent links. Thus, an appropriate form of the transformation function could be selected from the complete set. Based on the limit theorem, a new combination method, theoretical limit fusion (TLF) method, is proposed. We use the Parzen kernel method [16] of destiny estimation in the TLF method. Simulations and empirical studies have shown that TLF method can achieve higher prediction accuracy.

Section 2 introduces a mathematical description for the theoretical limit of combination methods and evaluation metrics for link prediction. Section 3 proposes and proves necessary and sufficient conditions for the theoretical limit of combination methods. Section 4 proposes a fusion link prediction method based on limit theorem (TLF method). Section 5 provides simulation examples for limit theorem and proposed TLF method with other combination methods, and gives comparison experiments in real networks. Sections 6 and 7 discuss some results and conclude the paper.

2. Problem Description and Evaluation Metrics

2.1. Problem Description

Given a network $G(V, E)$ at time t , where $V = \{v_1, v_2, \dots, v_N\}$ is the set of nodes and $E = \{e_1, e_2, \dots, e_M\}$ is the set of links. The observed links, E , are randomly divided into training, E^T , and probe, E^P , sets, where $E = E^T \cup E^P$ and $E^T \cap E^P = \emptyset$. Link prediction aims to predict missing links at current network or new links for a future time $t'(t' > t)$ [2]. Link prediction combination methods fuse several similarity indices and obtain a synthetic index and can be described in mathematic as follows. Let $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$ be the scores of existing links as given by n structural similarity indices, and follow probability density function (pdf) $f(x) = f(x_1, x_2, \dots, x_n)$. Let $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T$ be the scores of nonexistent links as n structural similarity indices, and follow $g(x) = g(x_1, x_2, \dots, x_n)$. We need to find the transformation function, $l(x)$, and obtain the synthetic score, $X = l(\mathbf{X})$, $Y = l(\mathbf{Y})$ that maximizes evaluation metrics. Figure 1 is the diagram of combination methods.

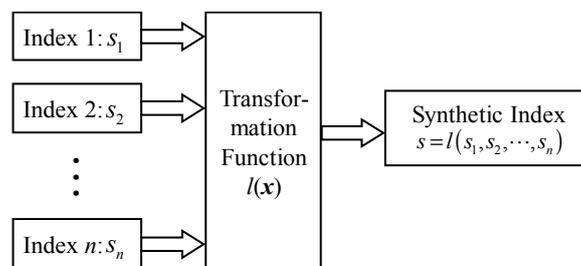


Figure 1. Combination methods.

2.2. Evaluation Metrics

Let the synthetic score $X = l(\mathbf{X})$ follow pdf $f_X(x)$, and $Y = l(\mathbf{Y})$ follow $g_Y(x)$. X and Y are independent. We have the following metrics.

2.2.1. Area under the Receiver Operation Characteristics Curve (AUC)

A receiver operating characteristics (ROC) curve is a two-dimensional depiction of classifier performance [17]. In the field of link prediction, the ROC curve abscissa represents the probability of nonexistent links i.e., the false positive rate (FPR), when the link prediction score is greater than some threshold, μ , and $FPR = \int_{\mu}^{\infty} g_Y(x)dx$. The ordinate represents the probability of missing links,

i.e., the true positive rate (TPR), when score $> \mu$, and $TPR = \int_{\mu}^{\infty} f_X(x)dx$, TPR is equivalent to Recall. According to [18], AUC can be derived as

$$\begin{aligned}
 P(X > Y) &= \iint_{X>Y} f_X(x)g_Y(y)dx dy \\
 &= \frac{1}{2} \iint_{X>Y} f_X(x)g_Y(y)dx dy + \frac{1}{2} \left(1 - \iint_{X \leq Y} f_X(x)g_Y(y)dx dy \right) \\
 &= \frac{1}{2} \iint \text{sgn}(x - y) f_X(x)g_Y(y)dx dy + \frac{1}{2} \\
 &= \frac{1}{2} \mathbb{E}[\text{sgn}(X - Y) + 1],
 \end{aligned}
 \tag{1}$$

where

$$\text{sgn}(x) = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases}
 \tag{2}$$

In the real network, original data is randomly divided into training set and the probe set. Equation (1) means that for n independent comparisons, if there are n' comparisons where the missing link returns a higher score and n'' comparisons where the missing and nonexistent links return the same score, we can obtain the algorithm expression of AUC:

$$\text{AUC} = \frac{n' + 0.5n''}{n}
 \tag{3}$$

2.2.2. Precision

Precision can be defined as the ratio of correct to (correct and error) prediction proportions when score $> \mu$, i.e.,

$$\begin{aligned}
 \text{Precision} &= \frac{P(\omega_1) \int_{\mu}^{+\infty} f_X(x)dx}{P(\omega_1) \int_{\mu}^{+\infty} f_X(x)dx + P(\omega_2) \int_{\mu}^{+\infty} g_Y(x)dx} \\
 &= \frac{P(\omega_1) \text{TPR}}{P(\omega_1) \text{TPR} + P(\omega_2) \text{FPR}}.
 \end{aligned}
 \tag{4}$$

In the real network, if the top L links are predicted ones, with m links being right (i.e., there are m links in E^P), then

$$\text{Precision} = \frac{m}{L}
 \tag{5}$$

Owing to the imbalance of positive and negative samples, link prediction usually uses AUC metric. In application, high Precision means target links are accurate, and these links can be used directly. AUC and Precision are two important metrics in link prediction, we will study the theoretical limit using the two metrics.

3. Theoretical Limit Theorem

Theorem 1. Let $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$ and $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T$ be random vectors following the joint distributions $f(\mathbf{x})$ and $g(\mathbf{x})$, respectively, where $m\{\mathbf{x} : f(\mathbf{x})/g(\mathbf{x}) = C, g(\mathbf{x}) \neq 0, \forall C \in \mathbb{R}\} = 0$. (m represents the measure of a set.) Then the following conditions are equivalent.

- (a) A monotonically increasing function $r(x)$ exists, such that $l(\mathbf{x}) = r[f(\mathbf{x})/g(\mathbf{x})]$, $g(\mathbf{x}) \neq 0$, a.e. $\mathbf{x} \in \mathbb{R}^n$.
- (b) Transformation function $l(\mathbf{x})$ produces maximum AUC. If we add a condition in Theorem that prior probability of existing and nonexistent links be $P(\omega_1)$ and $P(\omega_2)$, respectively. Then the following conditions are equivalent to (a) and (b):
- (c) for any α , there exists the corresponding threshold μ_1 for transformation $l(\mathbf{x})$, and satisfies $\alpha = P(\omega_1) \int_{\mu_1}^{+\infty} f_X(x)dx + P(\omega_2) \int_{\mu_1}^{+\infty} g_Y(x)dx$, such that transformation function $l(\mathbf{x})$ produces maximum Precision.

Proof. (a) \Rightarrow (b):

From the equivalent definition, AUC maximum is the maximum area under the ROC curve. For any FPR, if the TPRs corresponding to the ROC curve reach maximum, then the AUC reaches the maximum, i.e.,

$$\text{FPR} = \int_{\mu}^{\infty} g_Y(x)dx = \int_{E(l(x)>\mu)} g(x)dx, \tag{6}$$

$$\text{TPR} = \int_{\mu}^{+\infty} f_X(x)dx = \int_{E(l(x)>\mu)} f(x)dx, \tag{7}$$

where $E(l(x) > \mu)$ is a set $\{x \in \mathbb{R}^n : \mu \in \mathbb{R}, l(x) > \mu\}$, and $m\{x : l(x) = C, \forall C \in \mathbb{R}\} = 0$.

We use Lagrange’s undetermined multipliers to solve this problem. For any specified FPR (denoted as FPR_0), the TPR corresponding to the ROC curve reaches maximum is equivalent as φ reaches maximum,

$$\begin{aligned} \varphi &= \int_{E(l(x)>\mu)} f(x)dx + \lambda \left[\text{FPR}_0 - \int_{E(l(x)>\mu)} g(x)dx \right] \\ &= \lambda \text{FPR}_0 + \int_{E(l(x)>\mu)} [f(x) - \lambda g(x)]dx. \end{aligned} \tag{8}$$

Function φ will be maximized if we choose set $E(l(x) > \mu)$ such that the integrand is positive, i.e., if

$$f(x) - \lambda g(x) > 0, \tag{9}$$

then $x \in E(l(x) > \mu)$. Which means, no matter what is λ , if we select the set of x which makes the integrand $f(x) - \lambda g(x)$ always be positive, the function φ will reach maximum; if the set contains x that makes the integrand be negative, function φ will decrease. Let $l(x) = f(x)/g(x)$ and $\mu = \lambda$, and the set, $E(l(x) > \mu)$, equals to $E(f(x)/g(x) > \lambda)$, which satisfies (8), i.e.,

$$\varphi = \lambda \text{FPR}_0 + \int_{E(f(x)/g(x)>\lambda)} [f(x) - \lambda g(x)]dx \tag{10}$$

Thus, for any FPR, the TPR corresponding to the ROC curve reaches the maximum, so the AUC reaches the maximum when \mathbf{X} and \mathbf{Y} are transformed by $l(x) = f(x)/g(x)$.

Let $r(x)$ be a monotonically increasing function; and $h(x)$ be the inverse function of $r(x)$. If $h'(x) = 1/r'(x)$, then $h(x)$ and $r(x)$ have the same monotonicity, and both are increasing functions. Thus, $|h'(x)| = h'(x)$. The pdf of $X_2 = r(X_1)$ is $f_{X_2}(x) = f_{X_1}[h(x)]h'(x)$, and the pdf of $Y_2 = r(Y_1)$ is $g_{Y_2}(x) = g_{Y_1}[h(x)]h'(x)$. Thus,

$$\begin{aligned} \text{AUC} &= P(X_2 > Y_2) = \int_{-\infty}^{+\infty} f_{X_2}(x) \int_{-\infty}^x g_{Y_2}(y) dy dx \\ &= \int_{-\infty}^{+\infty} f_{X_1}(h(x))h'(x) \int_{-\infty}^x g_{Y_1}(h(y))h'(y) dy dx \\ &= \int_{-\infty}^{+\infty} f_{X_1}(x) \int_{-\infty}^x g_{Y_1}(y) dy dx \\ &= P(X_1 > Y_1). \end{aligned} \tag{11}$$

We have proved (a) \Rightarrow (b).

(b) \Rightarrow (a): If $l_2(x) \neq r[l(x)]$, where $r(x)$ is increasing function, there exists $l_2(x)$ such that \mathbf{X} , \mathbf{Y} transforming from $l_2(x)$ can also produce maximum AUC, and then the corresponding ROC curves are the same. Otherwise, if ROC curves are different, except the same part, for any FPR, there is at least a ROC curve which doesn’t reach maximum TPR, and contradict with maximum AUC. Since $m\{x : f(x)/g(x) = C, g(x) \neq 0, \forall C \in \mathbb{R}\} = 0$ and the ROC curve is the same for any point (FPR, TPR) on the two ROC curves, thus,

- i. For any $\text{FPR} \in [0, 1]$, and any μ_{FPR} , there exist $\mu_{2\text{FPR}}$, such that $E(l(x) > \mu_{\text{FPR}}) = E(l_2(x) > \mu_{2\text{FPR}})$ for a.e. $x \in \mathbb{R}^n$;

- ii. For any $\mu_{FPR}^* > \mu_{FPR}$, if $E(l(x) > \mu_{FPR}^*) = E(l_2(x) > \mu_{2FPR}^*)$ and $E(l(x) > \mu_{FPR}) = E(l_2(x) > \mu_{2FPR})$, then $\mu_{2FPR}^* > \mu_{2FPR}$.

Let $y_1 = l(x)$, then a set of y_1 exist with nonzero measure, such that $l_2(x) \neq r[l(x)]$, i.e., $m\{y_1 : l_2(x) \neq r[l(x)]\} \neq 0$. Let $\sigma = \{y_1 : l_2(x) \neq r[l(x)]\}$. If $y_1 \in \sigma$, $l_2(x), l_1(x)$ satisfies function relation $l_2(x) = s[l(x)]$, but $s(x)$ is not increasing, then for any $\mu_1 \in \sigma$, condition (ii) does not hold. If $y_1 \in \sigma$, $l_2(x)$ and $l(x)$ are not functionally related, then neither condition (i) or (ii) hold. Thus (b) \Rightarrow (a) is established.

(c) \Leftrightarrow (b): Let $k = \text{TPR}/\text{FPR}$ be the slope of the secant for any point on the ROC curve to the origin, then Precision = $k/(k + \lambda)$, $\lambda = P(\omega_2)/P(\omega_1)$. For any α , that $l(x)$ produces maximum Precision is equivalent that k reaches maximum. And equivalent that for any α , $\alpha = P(\omega_1) \int_{\mu_1}^{+\infty} f_X(x)dx + P(\omega_2) \int_{\mu_1}^{+\infty} g_Y(x)dx$, $\text{TPR} = \int_{\mu_1}^{+\infty} f_X(x)dx$ is maximum. Since this condition is established for any α , then it is equivalent that for any $\text{FPR} \in [0, 1]$, the corresponding TPR reaches maximum, and equivalent to $l(x)$ produces maximum AUC. \square

Note 1: the condition $m\{x : f(x)/g(x) = C, g(x) \neq 0, \forall C \in \mathbb{R}\} = 0$ is for exclusion that when $f(x)/g(x) = C$, (C is a constant), transformation function can be defined randomly on set $\sigma = \{x : f(x)/g(x) = C, g(x) \neq 0, \forall C \in \mathbb{R}\} \cap \mathbb{R}^n$. For example, let us construct the pdf of some random vector \tilde{X} as

$$\tilde{f}(x) = \begin{cases} f(x), & x \in \mathbb{R}^n \setminus \sigma \\ kg(x), & x \in \sigma \end{cases}, k \in \mathbb{R}, k < \frac{f(x)}{g(x)}, \tag{12}$$

$\int_{\mathbb{R}^n} \tilde{f}(x)dx = 1$. Let the transformation function be

$$l(x) = \begin{cases} \frac{f(x)}{g(x)}, & x \in \mathbb{R}^n \setminus \sigma \\ l^*(x), & x \in \sigma \end{cases} \tag{13}$$

then no matter how $l^*(x)$ is defined, only when $l^*(x) < \min[f(x)/g(x)]$ can the $l(x)$ produce the maximum AUC of (\tilde{X}, Y) . In particular, if $f(x) = g(x), x \in \mathbb{R}^n$, regardless of how $l(x)$ is defined, $\text{AUC} = 0.5$. Thus, maximum $\text{AUC} = 0.5$.

Note 2: The arbitrariness of the ratio α must be emphasized in condition (c). If we omitted “any α ”, then (b) \Rightarrow (c) can be established but (c) \Rightarrow (b) cannot. The meaning of α in application is a ratio of the whole data, for any $l(x)$, a ratio α corresponds a threshold μ .

Theorem shows that no matter which evaluation criteria choose, transformation functions that provide maximum link prediction accuracy constitute a function cluster, $\Phi = \{l(x) : l(x) = r[f(x)/g(x)], g(x) \neq 0\}$, where $r(x)$ is a monotonically increasing function. Therefore, the accuracy of the combination method must be greater than or equal to the accuracy of each single dimension.

4. A Fusion Link Prediction Method Based on Limit Theorem

4.1. The Algorithm

The limit theorem of combination method shows that when selecting transformation function as $l(x) = f(x)/g(x)$ or its monotone increasing transformation, the AUC and Precision of synthetic score reaches the maximum. In the real network, because $f(x)$ and $g(x)$ are unknown, the pdfs need to be estimated from multidimensional data. Let the estimated pdfs be $\hat{f}(x)$ and $\hat{g}(x)$. On the basis of limit theorem, we define the transformation function as the ratio of estimated pdfs, i.e.,

$$\hat{l}(x) = \hat{f}(x)/\hat{g}(x) \tag{14}$$

then we obtained the synthetic score, $s = \hat{l}(x)$, and used for link prediction. This method is called theoretical limit fusion (TLF) method.

Before evaluating $f(\mathbf{x})$ and $g(\mathbf{x})$, the input link prediction scores need to be normalized,

$$s_k^*(i, j) = \frac{0.5N^2 \cdot s_k(i, j)}{\sum_{i=1}^N \sum_{j=1}^N s_k(i, j)}, k = 1, 2, \dots, d \tag{15}$$

$s_k(i, j)$ represents the k -th similarity score for node pair i, j . N is the dimension of adjacent matrix, and d is the number of similarity indices.

The limit theorem of combination method transformed the link prediction indices fusion problem into the estimation of pdfs. Statistical methods for estimating density functions can be applied to this problem, directly. The Parzen kernel method [16] of destiny estimation is used in this paper. The multivariate kernel density estimate defined as:

$$\hat{f}(\mathbf{x}) = \frac{1}{n_s h^d} \sum_{i=1}^{n_s} K\left[\frac{1}{h}(\mathbf{x} - \mathbf{x}_i)\right] \tag{16}$$

where h is the window width, n_s is the sample size, and $K(\mathbf{x})$ is a multivariate kernel defined for d -dimensional \mathbf{x} , such that

$$\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1 \tag{17}$$

A form of the pdf estimate commonly used is Gauss kernel,

$$K(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \exp\left(\frac{-\mathbf{x}^T \mathbf{x}}{2}\right) \tag{18}$$

In summary, the steps of TLF are listed as Table 1.

Table 1. The steps of theoretical limit fusion (TLF) method.

Step 1	Divide the network into training set, E^T , and probe set, E^P ;
Step 2	According to Equation (15), normalize these similarity indices, then we distinguish existing links and nonexistent links in the training set;
Step 3	Based on Equation (16), estimate the pdfs of existing links and nonexistent links, and we obtain the estimated pdfs as $\hat{f}(\mathbf{x})$ and $\hat{g}(\mathbf{x})$;
Step 4	Obtain the synthetic score of n structure similarity indices according to Equation (14);
Step 5	Calculate the accuracy such as AUC metric or Precision metric on the probe set.

4.2. Complexity Analysis

For a given undirected, unweighted graph $G(V, E)$, let $N = |V|$ be the number of nodes and let $m = |E|$ be the number of edges, and let n_s be the sample size. During the estimation of pdfs in (16), the entire samples are scanned once. A scan of samples requires time $O(d \cdot n_s)$ and it is less than $O(N^2)$. This is the step of model training or pdf estimation. Among all combination methods, there is an inevitable time complexity, that is to obtain the similarity matrix or final link prediction scores according to Equation (14). This step requires time $O(d \cdot n_s \cdot N^2)$. So, the TLF method will take time more than $O(N^3)$. The main space needs to storage estimator and adjacent matrix or final similarity matrix. The spatial complexity is $O(N^2)$.

5. Simulation and Experiment

We programmed the algorithm using Matlab (MathWorks, Beijing, China, 2014), and runs on a single machine equipped with RedHat 6.4. The host memory is 16 GB, with 3.4 GHz CPU, and the Matlab version is R2014b. In simulations from Section 5.1, 4-dimensional pdfs are supported to verify limit theorem and the effectiveness of TLF method. We also test the resulting method in real networks. We use TLF method to fuse 4 local similarity indices, CN [19], AA [20], RA and PA [21,22]. These indices are 4 simple indices with low computation complexity about $O(N \cdot \langle k \rangle^2)$, where $\langle k \rangle$ represents the average degree of nodes in a network. CN index only considers common neighbors of node pairs; PA index only considers the degree of two nodes; AA and RA consider both common neighbors and degree of nodes with different weights. And compare the method with fusion methods such as naïve Bayes and logistic regression and other global indices with computation complexity more than $O(N^3)$.

5.1. Simulation Examples

Four types of structural similarity indices were simulated to evaluate node pairs with and without links. The pdfs of the structural similarity indices are also provided. We construct 3 groups of known distributions for the similarity indices pdfs. One thousand samples extracted from 10,000 existing links and 100,000 samples of nonexistent links were generated following the appropriate pdfs. The 1000 samples serve as probe set; the 100,000 samples with 1000 probe links serve as unknown links for training; and the remaining 9000 samples serve as train set of existing links. Each sample had 4 dimensions to simulate 4 similarity scores. We first compute AUC and Precision for each dimension, then use proposed TLF method to obtain the synthetic score and calculate the AUC and Precision, compared with other combination methods such as Naïve Bayes and logistic regression. Finally, we calculate AUC and Precision using the theoretical limit theorem and compare with the above methods.

Let random vectors $\mathbf{X} = (X_1, X_2, X_3, X_4)^T$ and $\mathbf{Y} = (Y_1, Y_2, Y_3, Y_4)^T$ be the scores of existing and nonexistent links, which follow $f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$ and $g(\mathbf{x}) = g(x_1, x_2, \dots, x_n)$ pdfs, respectively.

Let $f(\mathbf{x}), g(\mathbf{x})$ are 4-dimensional normal distributions,

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}, \tag{19}$$

where $\text{diag}(\Sigma)\mathbf{1} = (\sigma_1^2, \sigma_2^2, \sigma_3^2, \sigma_4^2)^T$, and $\Sigma_{ij} = r_{ij}\sigma_i\sigma_j$.

The parameter sets for the 2 groups of simulation examples are as follows.

Group 1: $\Theta_{1f} = \{\boldsymbol{\mu}_{1f}, \Sigma_{1f}\}$, and $\Theta_{1g} = \{\boldsymbol{\mu}_{1g}, \Sigma_{1g}\}$;

Group 2: $\Theta_{2f} = \{\boldsymbol{\mu}_{2f}, \Sigma_{2f}\}$ and $\Theta_{2g} = \{\boldsymbol{\mu}_{2g}, \Sigma_{2g}\}$.

In each group, $\boldsymbol{\mu}_{1f} = (1, 2, 1.7, 2.1)^T$, $\boldsymbol{\mu}_{1g} = (1.3, 2.5, 2.1, 2.8)^T$, $\boldsymbol{\mu}_{2f} = (1, 2, 1.7, 2.1)^T$, $\boldsymbol{\mu}_{2g} = (1.5, 3.5, 2.8, 3)^T$, $\text{diag}(\Sigma_{1f})\mathbf{1} = (1.5^2, 2.2^2, 3^2, 2.5^2)^T$, $\text{diag}(\Sigma_{1g})\mathbf{1} = (2^2, 2.2^2, 3^2, 2.5^2)^T$, $\text{diag}(\Sigma_{2f})\mathbf{1} = (1.5^2, 2.2^2, 3^2, 2.5^2)^T$, $\text{diag}(\Sigma_{2g})\mathbf{1} = (2.5^2, 3.5^2, 4^2, 2.5^2)^T$, $r_{1f} = r_{1g} =$

$$\begin{bmatrix} 1 & 0.8 & 0.76 & 0.56 \\ 0.8 & 1 & 0.85 & 0.74 \\ 0.76 & 0.85 & 1 & 0.93 \\ 0.56 & 0.74 & 0.93 & 1 \end{bmatrix}, \text{ and } r_{2f} = r_{2g} = \begin{bmatrix} 1 & 0.62 & 0.45 & 0.34 \\ 0.62 & 1 & 0.28 & 0.47 \\ 0.45 & 0.28 & 1 & 0.65 \\ 0.34 & 0.47 & 0.65 & 1 \end{bmatrix}.$$

The window width h of TLF method in the group 1 and 2 is $h = 0.1$.

Group 3: Let

$$f_3(x) = x_1x_2x_3x_4 + x_1x_4 + x_3 \exp(x_1) \log(x_2) \quad (20)$$

$$(0 \leq x_1 \leq 3, 1 \leq x_2 \leq 3, 3 \leq x_3 \leq 5, 2 \leq x_4 \leq 3.5)$$

and

$$g_3(x) = x_1x_2x_3x_4 + x_3 \exp(x_1) \log(x_2) \quad (21)$$

$$(0 \leq x_1 \leq 4, 1 \leq x_2 \leq 3, 3 \leq x_3 \leq 5, 2.5 \leq x_4 \leq 5)$$

We ignore the constant that makes the integral of $f(x)$, $g(x)$ equal to 1. The simulation results of group 3 are shown as Table 2.

Table 2. Simulation results of group 1 and group 2. The bold figure indicates the best accuracy in each dimension and combination method.

Parameters	Accuracy	Dim1	Dim2	Dim3	Dim4	NB	LR	TLF	Theoretical Limit	Transform by Increasing Function
Group 1	AUC	0.554	0.566	0.547	0.585	0.610	0.668	0.691	0.738	0.738
	Precision	0.047	0.015	0.014	0.027	0.038	0.020	0.097	0.120	0.120
Group 2	AUC	0.569	0.660	0.604	0.622	0.765	0.676	0.786	0.792	0.792
	Precision	0.114	0.140	0.081	0.038	0.153	0.051	0.212	0.241	0.241

The window width h of TLF method in the group 3 is $h = 0.1$.

The simulation results in Tables 2 and 3 show us that we can calculate the theoretical limit of combination method based on Theorem 1, and the limit AUC and Precision are highest among all listed methods, though we cannot list all possible conditions. Results also show that TLF method can fuse the information effectively, and obtain the optimum accuracy. We also verify that the transformation of monotonically increasing function does not change the theoretical limit. Theorem 1 provides a platform that can compare each combination method by constructing some distributions, and direct an effect combination method TLF.

Table 3. Simulation results of group 3. The bold figure indicates the best accuracy in each dimension and combination method.

Accuracy	Dim1	Dim2	Dim3	Dim4	NB	LR	TLF	Theoretical Limit	Transform by Increasing Function
AUC	0.770	0.505	0.488	0.878	0.938	0.923	0.950	0.956	0.956
Precision	0.567	0.007	0.007	0.654	0.711	0.100	0.815	0.858	0.858

5.2. Experiments in Real Networks

The significance of simulation is that the theoretical limit can be derived by theoretical calculation or numerical calculation, and all combination methods can be used to compare with it, finding shortcomings and gaps to design a more rational method. However, the simulation data is different from real network data. We use TLF method to fuse several similarity indices and test in real networks. The basic similarity indices we use are Common Neighbor index (CN) [19], Adamic-Adar index (AA) [15], Resource Allocation index (RA) and Preferential Attachment index (PA) [21,22]. These indices are local indices. Several global indices such as Katz index [23], Average Commute Time index (ACT) and Cosine Similarity Time index (Cos+) are served as comparisons [24,25]. The definitions of the above indices and their meanings are listed as Table 4.

Table 4. Definitions and descriptions of similarity indices.

Index	Equation	Description
CN	$s_{CN}(i, j) = \Gamma(i) \cap \Gamma(j) $	$\Gamma(i)$ is the set of neighbors of node i . $ \cdot $ represents cardinality of a set. CN index denotes the common neighbors between nodes i and j .
AA	$s_{AA}(i, j) = \sum_{z \in \Gamma(i) \cap \Gamma(j)} \frac{1}{\log k_z}$	AA index weights the common neighbors by the reciprocal of the logarithm of each node's degree.
RA	$s_{RA}(i, j) = \sum_{z \in \Gamma(i) \cap \Gamma(j)} \frac{1}{k_z}$	RA index weights the common neighbors by the reciprocal of each node's degree.
PA	$s_{PA}(i, j) = k_i k_j$	PA index expresses preferential attachment by node's degree.
Katz	$s_{Katz}(i, j) = \left[\lim_{n \rightarrow \infty} \sum_{m=1}^n (\alpha \mathbf{A})^m \right]$	\mathbf{A} is adjacent matrix of network. Katz index considers all path between two nodes and gives more weights, α , to the shorter paths.
ACT	$s_{ACT}(i, j) = \frac{1}{l_{ii}^+ + l_{jj}^+ - 2l_{ij}^+}$	l_{xy}^+ is the corresponding element in \mathbf{L}^+ , and \mathbf{L}^+ denotes the pseudo-inverse of laplacian matrix.
Cos+	$s_{Cos+}(i, j) = \frac{v_i^T v_j}{ v_i \cdot v_j } = \frac{l_{ij}^+}{\sqrt{l_{ii}^+ \cdot l_{jj}^+}}$	According to \mathbf{L}^+ , Cos+ calculates cosine similarity of two vectors in matrix \mathbf{L}^+ .

We use TLF method to fuse 4 local similarity indices, and compare with fusion method such as naïve Bayes and logistic regression and other global indices. Our experiments are performed on 11 different real networks. (1) Food Web Everglades Web (FWEW) [26]; (2) Food Web Florida Bay(FWFB) [27]; (3) Protein-protein Interactions Cell (PPI Cell) [28]; (4) CKM-3 [29]; (5) Netscience (NS) [30]; (6) Yeast [31]; (7) Political Blogosphere(PB) [32]; (8) Email [33]; (9) CA-GrQc(CG) [34]; (10) Com-dblp(CD) [35]; (11) Email Enron (EE) [36,37]. The basic topological features of 11 real networks are listed in Table 5. Each original data is randomly divided into training set of 90% links, and the probe set of 10% links.

Tables 6 and 7 show the comparisons between TLF method and other combination methods or global indices using AUC and Precision metrics. Each result is the average of 10 realizations. When calculating the Precision metric in Equation (5), we take $L = 100$ in datasets 1 to 8, and take $L = 1000$ in datasets 9 to 11. In the large networks, TLF method needs to sample to save the computing sources, and in datasets 10 to 11, the under-sampling rate is set as 1000.

Table 5. Basic topological features of 6 example networks. $|V|$ and $|E|$ are the total numbers of nodes and links, respectively. $\langle k \rangle$ represents the average degree of nodes in a network. C and r are the clustering coefficient and assortative coefficient respectively. H is the degree heterogeneity, defined as $H = \frac{\langle k^2 \rangle}{\langle k \rangle^2}$.

Data	$ V $	$ E $	$\langle k \rangle$	r	C	H
FWEW	69	880	25.51	-0.298	0.560	1.275
FWFB	128	2075	32.42	-0.112	0.335	1.24
PPI_Cell	127	237	3.732	0.035	0.455	1.649
CKM-3	246	423	3.439	0.102	0.356	1.335
Yeast	2375	11,693	9.85	0.469	0.378	3.48
PB	1222	16,717	27.36	-0.221	0.361	2.97
NS	1589	2742	3.451	0.462	0.889	2.011
Email	1133	5451	9.622	0.078	0.297	1.942
CG	5242	14,496	1.11	0.659	0.720	3.05
CD	425,957	1,049,866	4.93	0.267	0.267	4.412
EE	36,692	183,831	10.02	-0.111	0.746	13.9

Table 6. Comparisons of the AUC value between TLF and other combination methods or global indices. In each network, the selected window width h is along with the AUC value. The bold figure indicates the best AUC.

Data	CN	AA	RA	PA	ACT	Cos+	Katz	NB	LR	TLF
FWEW	0.687	0.694	0.714	0.819	0.793	0.511	0.727	0.825	0.832	0.876 ($h = 0.1$)
FWFB	0.624	0.624	0.624	0.742	0.727	0.649	0.680	0.749	0.762	0.781 ($h = 0.1$)
PPI_Cell	0.736	0.745	0.740	0.699	0.779	0.783	0.822	0.753	0.679	0.831 ($h = 0.3$)
CKM-3	0.661	0.665	0.661	0.585	0.560	0.535	0.928	0.683	0.675	0.713 ($h = 0.15$)
Yeast	0.918	0.918	0.915	0.869	0.903	0.958	0.962	0.925	0.934	0.968 ($h = 0.2$)
PB	0.922	0.928	0.928	0.906	0.890	0.932	0.934	0.931	0.936	0.949 ($h = 0.3$)
NS	0.994	0.994	0.995	0.709	0.558	0.507	0.996	0.998	0.999	0.999 ($h = 0.2$)
Email	0.849	0.852	0.851	0.817	0.801	0.889	0.908	0.865	0.870	0.912 ($h = 0.15$)
CG	0.966	0.965	0.967	0.992	0.549	0.679	0.996	0.984	0.991	0.994 ($h = 0.1$)
CD	0.962	0.968	0.971	0.943	0.912	0.971	0.915	0.975	0.973	0.982 ($h = 0.15$)
EE	0.981	0.984	0.984	0.927	0.903	0.980	0.514	0.985	0.987	0.992 ($h = 0.15$)

Table 7. Comparisons of the Precision value between TLF and other combination methods or global indices. In each network, the corresponding window width h is the same as Table 6. The bold figure indicates the best Precision.

Data	CN	AA	RA	PA	ACT	Cos+	Katz	NB	LR	TLF
FWEW	0.143	0.145	0.162	0.334	0.271	0.004	0.196	0.301	0.325	0.543
FWFB	0.071	0.072	0.083	0.240	0.184	0.029	0.148	0.249	0.283	0.382
PPI_Cell	0.052	0.048	0.073	0.012	0.045	0.061	0.058	0.072	0.068	0.085
CKM-3	0.051	0.059	0.062	0.011	0.001	0.003	0.061	0.060	0.062	0.064
Yeast	0.652	0.703	0.461	0.439	0.487	0.291	0.721	0.712	0.723	0.785
PB	0.381	0.320	0.212	0.100	0.129	0.298	0.381	0.411	0.395	0.452
NS	0.820	0.971	0.982	0.008	0.004	0.006	0.823	0.988	0.986	0.991
Email	0.202	0.253	0.214	0.039	0.031	0.086	0.231	0.263	0.289	0.347
CG	0.972	0.969	0.967	0.991	0.557	0.663	0.998	0.983	0.989	0.996
CD	0.901	0.924	0.931	0.892	0.867	0.937	0.912	0.939	0.942	0.951
EE	0.981	0.984	0.987	0.924	0.898	0.912	0.516	0.988	0.985	0.992

The results show us that TLF method performs better than other fusion methods such as naïve Bayes and logistic regression, no matter what evaluation metric use. Almost all combination methods are better than 4 basic indices. From the limit theorem, combination methods are dependent with each dimension. The promotion of fusion index is restrict to each similarity index. Experiment results also exposed this problem: if the single similarity indices perform not well, the fusion index cannot significantly improve the accuracy. For example, in the CKM-3 network, though we use TLF method to fuse 4 basic similarity indices can improve the AUC obviously, it cannot be better than Katz index (0.928).

6. Discussion

Many combination methods try to find the nonlinear relation of every dimensions, and want to obtain a more reasonable fusion function to promote the prediction accuracy. For example, link prediction method based on the choquet fuzzy integral [7] uses fuzzy measures to measure the importance of each similarity index in the fusion process and the interaction between them. Logistic regression based index adopts logistic function to learn the relation of multiple structural features and obtain an adaptive link prediction method [38]. In fact, according to the limit theorem, the nonlinear relation is the ratio of two joint probability destiny functions or its monotone increasing transformation. The best fusion function is a measurement of difference between existing and nonexistent links, and it reflects the relativity of existing and nonexistent links. The essence of combination methods is trying to approximate the pdfs from many aspects. Limit theorem provides a unified interpretation for all combination methods. On the basis of theoretical limit theorem,

the proposed TLF method evaluates two pdfs directly, and it has a better fusion effect from results of simulation and experiment in real network.

7. Conclusions

This paper proposes mathematic description of link prediction combination methods and derives the limit theorem. Before the mathematic description we proposed, many combination methods have been put forward and widely used. However, all these methods are groping respectively without unified explanation. Limit theorem solved this problem and provided a guidance for link prediction method design. The TLF method based on limit theorem can achieve higher prediction accuracy.

Acknowledgments: We acknowledge professor Guo'en Hu for inspirations. This work was partially supported by the Foundation for Innovative Research Groups of the National Natural Science Foundation of China (No. 61521003), and National Natural Science Foundation of China (No. 61601513).

Author Contributions: Yiteng Wu and Hongtao Yu proposed mathematical description of combination method; Yiteng Wu proposed and proved the theoretical limit theorem; Yiteng Wu and Ruiyang Huang designed the experiments and analyzed the results. Yingle Li and Senjie Lin wrote part of code.

Conflicts of Interest: The authors declare no conflict of interest. The founding sponsors had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and in the decision to publish the results.

References

- Seife, C. What are the limits of conventional computing. *Science* **2005**, *309*, 96. [CrossRef] [PubMed]
- Wang, P.; Xu, B.; Wu, Y.; Zhou, X. Link prediction in social networks: The state-of-the-art. *Sci. China Inf. Sci.* **2015**, *58*, 1–38. [CrossRef]
- Lü, L.; Pan, L.; Zhou, T.; Zhang, Y.-C.; Stanley, H.E. Toward link predictability of complex networks. *Proc. Natl. Acad. Sci.* **2015**, *112*, 2325–2330. [CrossRef] [PubMed]
- Lü, L.; Zhou, T. Link prediction in complex networks: A survey. *Phys. A Stat. Mech. Appl.* **2011**, *390*, 1150–1170. [CrossRef]
- Wohlfarth, T.; Ichise, R. Semantic and Event-Based Approach for Link Prediction. In Proceedings of the Practical Aspects of Knowledge Management (PAKM 2008), Yokohama, Japan, 22–23 November 2008. [CrossRef]
- Chiancone, A.; Franzoni, V.; Li, Y.; Markov, K.; Milani, A. Leveraging Zero Tail in Neighbourhood for Link Prediction. In Proceedings of the 2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT), Singapore, 6–9 December 2015; pp. 135–139. [CrossRef]
- Yu, H.T.; Wang, S.H.; Ma, Q. Link prediction algorithm based on the Choquet fuzzy integral. *Intell. Data Anal.* **2016**, *20*, 809–824. [CrossRef]
- He, Y.-l.; Liu, J.N.K.; Hu, Y.-X.; Wang, X.-Z. OWA operator based link prediction ensemble for social network. *Expert Syst. Appl.* **2015**, *42*, 21–50. [CrossRef]
- Liao, L.; He, X.; Zhang, H.; Chua, T.-S. Attributed Social Network Embedding. *Trans. Knowl. Data Eng.* **2017**. Available online: <http://www.comp.nus.edu.sg/~xiangan/papers/attributed-social-network-embedding.pdf> (accessed on 5 September 2017).
- Grover, A.; Leskovec, J. node2vec: Scalable Feature Learning for Networks. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, 13–17 August 2016.
- Wang, Z.; Chen, C.; Li, W. Predictive Network Representation Learning for Link Prediction. In Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, 7–11 August 2017; pp. 969–972.
- Chuan, P.M.; Son, L.H.; Ali, M.; Khang, T.D.; Huong, L.T.; Dey, N. Link prediction in co-authorship networks based on hybrid content similarity metric. *Appl. Intell.* **2017**. [CrossRef]
- Franzoni, V.; Chiancone, A.; Milani, A. A Multistrain Bacterial Diffusion Model for Link Prediction. *Int. J. Pattern Recognit. Artif. Intell.* **2017**, *31*, 1759024. [CrossRef]

14. Liu, B.; Sun, C.; Liu, M.; Wang, X.; Liu, F. Deep Belief Network-Based Approaches for Link Prediction in Signed Social Networks. *Entropy* **2015**, *17*, 2140–2169. [[CrossRef](#)]
15. Ou, Q.; Jin, Y.D.; Zhou, T.; Wang, B.H.; Yin, B.Q. Power-law strength-degree correlation from resource-allocation dynamics on weighted networks. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* **2007**, *75*, 021102. [[CrossRef](#)] [[PubMed](#)]
16. Parzen, E. On estimation of a probability density function and mode. *Ann. Math. Stat.* **1962**, *33*, 1065–1076. [[CrossRef](#)]
17. Fawcett, T. An introduction to ROC analysis. *Pattern Recognit. Lett.* **2006**, *27*, 861–874. [[CrossRef](#)]
18. Hanley, J.A.; McNeil, B.J. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* **1982**, *143*, 29–36. [[CrossRef](#)] [[PubMed](#)]
19. Lorrain, F.; White, H.C. Structural equivalence of individuals in social networks. *Soc. Netw.* **1977**, *1*, 67–98. [[CrossRef](#)]
20. Adamic, L.A.; Adar, E. Friends and neighbors on the web. *Soc. Netw.* **2003**, *25*, 211–230. [[CrossRef](#)]
21. Zhou, T.; Lü, L.; Zhang, Y.C. Predicting missing links via local information. *Eur. Phys. J. B-Condens. Matter Complex Syst.* **2009**, *71*, 623–630. [[CrossRef](#)]
22. Barabasi, A.L.; Albert, R. Emergence of scaling in random networks. *Science* **1999**, *286*, 509–512. [[CrossRef](#)] [[PubMed](#)]
23. Coleman, J.; Katz, E.; Menzel, H. The Diffusion of an Innovation among Physicians. *Sociometry* **1957**, *20*, 253–270. [[CrossRef](#)]
24. Klein, D.J.; Randić, M. Resistance distance. *J. Math. Chem.* **1993**, *12*, 81–95. [[CrossRef](#)]
25. Fouss, F.; Pirotte, A.; Renders, J.-M.; Saerens, M. Random-Walk Computation of Similarities between Nodes of a Graph with Application to Collaborative Recommendation. *IEEE Trans. Knowl. Data Eng.* **2007**, *19*, 355–369. [[CrossRef](#)]
26. Ulanowicz, R.E.; DeAngelis, D.L.; Egnotovitch, M.S. Network Analysis of Trophic Dynamics in South Florida Ecosystems, FY 99: The Graminoid Ecosystem. 2000. Available online: https://www.researchgate.net/publication/237005295_Network_Analysis_of_Trophic_Dynamics_in_South_Florida_Ecosystems_FY_99_The_Graminoid_Ecosystem (accessed on 13 June 2016).
27. Ulanowicz, R.E.; Bondavalli, C.; Egnotovitch, M.S. *Network Analysis of Trophic Dynamics in South Florida Ecosystem, FY 97: The Florida Bay Ecosystem*; Technical Report; CBL: Chattanooga, TN, USA, 1998; pp. 98–123.
28. Kolaczyk, E.D. *Statistical Analysis of Network Data: Methods and Models*; Springer: New York, NY, USA, 2009. [[CrossRef](#)]
29. Coleman, J.; Katz, E.; Menzel, H. The Diffusion of an Innovation among Physicians 1. *Soc. Netw.* **1977**, *20*, 107–124. [[CrossRef](#)]
30. Newman, M.E.J. Finding community structure in networks using the eigenvectors of matrices. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* **2006**, *74* (Pt 2), 036104. [[CrossRef](#)] [[PubMed](#)]
31. Von Mering, C.; Krause, R.; Snel, B.; Cornell, M.; Oliver, S.G.; Fields, S.; Bork, P. Comparative assessment of large-scale data sets of protein protein interactions. *Nature* **2002**, *417*, 399–403. [[CrossRef](#)] [[PubMed](#)]
32. Adamic, L.A.; Glance, N. The political blogosphere and the 2004 U.S. election: Divided they blog. In Proceedings of the 3rd International Workshop on Link Discovery, Chicago, IL, USA, 21–25 August 2005; ACM: New York, NY, USA, 2005; pp. 36–43. [[CrossRef](#)]
33. Michalski, R.; Palus, S.; Kazienko, P. Matching Organizational Structure and Social Network Extracted from Email Communication. In *Business Information Systems*; Springer: Berlin, Germany, 2011; pp. 197–206.
34. Leskovec, J.; Kleinberg, J.; Faloutsos, C. Graph Evolution: Densification and Shrinking Diameters. *ACM Trans. Knowl. Discov. Data ACM TKDD* **2007**, *1*. [[CrossRef](#)]
35. Yang, J.; Leskovec, J. Defining and Evaluating Network Communities based on Ground-truth. In Proceedings of the 12th International Conference on Data Mining (ICDM), Brussels, Belgium, 10–13 December 2012. [[CrossRef](#)]
36. Leskovec, J.; Lang, K.J.; Dasgupta, A.; Mahoney, M.W. Community Structure in Large Networks: Natural Cluster Sizes and the Absence of Large Well-Defined Clusters. *Internet Math.* **2009**, *6*, 29–123. [[CrossRef](#)]

37. Klimmt, B.; Yang, Y. Introducing the Enron corpus. In Proceedings of the CEAS Conference 2004, Mountain View, CA, USA, 30–31 July 2004.
38. Ma, C.; Bao, Z.K.; Zhang, H.F. Improving link prediction in complex networks by adaptively exploiting multiple structural features of networks. *Phys. Lett. A* **2017**, *381*, 3369–3376. [[CrossRef](#)]



© 2017 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).