

Article

A novel probability model for LncRNA-Disease Association Prediction based on the Naïve Bayesian Classifier

Jingwen Yu ¹, Pengyao Ping ¹, Lei Wang ^{1,2*} , Linai Kuang ^{1,2}, Xueyong Li ² and Zhelun Wu ³

¹ Key Laboratory of Intelligent Computing & Information Processing, Xiangtan University, XiangTan 411105, People's Republic of China; jingwen-yu@qq.com; pengyao.ping@qq.com; wanglei@xtu.edu.cn; kla@xtu.edu.cn

² College of Computer Engineering & Applied Mathematics, Changsha University, Changsha 410001, Hunan, People's Republic of China; wanglei@xtu.edu.cn; kla@xtu.edu.cn; 137493255@qq.com

³ Department of Computer Science, Princeton University, Princeton, New Jersey, USA; zhelunw@princeton.edu

* Correspondence: wanglei@xtu.edu.cn; Tel.: +86-151-1110-9999

Academic Editor: name

Version June 24, 2018 submitted to Journal Not Specified

1. Method

1.1. Method for Applying the Naïve Bayesian Theory into GN_1

Based on the naïve Bayesian classifier, the posterior probabilities for an edge $e_{l_i-d_j}$, representing whether the node l_i is connected to d_j in GN_1 , are defined as follows:

$$\begin{aligned} p(e_{l_i-d_j} = 1 | CN(l_i, d_j)) &= \frac{p(m_1, m_2, \dots, m_h | e_{l_i-d_j} = 1) p(e_{l_i-d_j} = 1)}{p(CN(l_i, d_j))} \\ &= \frac{p(e_{l_i-d_j} = 1)}{p(CN(l_i, d_j))} \prod_{m_\delta \in CN(l_i, d_j)} p(m_\delta | e_{l_i-d_j} = 1) \end{aligned} \quad (1)$$

$$\begin{aligned} p(e_{l_i-d_j} = 0 | CN(l_i, d_j)) &= \frac{p(m_1, m_2, \dots, m_h | e_{l_i-d_j} = 0) p(e_{l_i-d_j} = 0)}{p(CN(l_i, d_j))} \\ &= \frac{p(e_{l_i-d_j} = 0)}{p(CN(l_i, d_j))} \prod_{m_\delta \in CN(l_i, d_j)} p(m_\delta | e_{l_i-d_j} = 0). \end{aligned} \quad (2)$$

From Equations (1) and (2), we can directly identify whether an lncRNA node is connected with a disease node or not in GN_1 . However, since it is often too complicated to calculate the value of $p(CN(l_i, d_j))$, we first define the probability of a potential association existing between l_i and d_j in GN_1 as follows:

$$S1(l_i, d_j) = \frac{p(e_{l_i-d_j} = 1 | CN(l_i, d_j))}{p(e_{l_i-d_j} = 0 | CN(l_i, d_j))} = \frac{p(e_{l_i-d_j} = 1)}{p(e_{l_i-d_j} = 0)} \prod_{m_\delta \in CN(l_i, d_j)} \frac{p(m_\delta | e_{l_i-d_j} = 1)}{p(m_\delta | e_{l_i-d_j} = 0)}, \quad (3)$$

where $p(m_\delta | e_{l_i-d_j}=1)$ and $p(m_\delta | e_{l_i-d_j}=0)$ are the conditional probabilities of a node m_δ belonging to $CN(l_i, d_j)$; they represent the possibilities of whether the node is a common neighboring node between

l_i and d_j in GN_1 or not, respectively. Moreover, according to Bayesian theory, these two conditional probabilities can be expressed as:

$$p(m_\delta | e_{l_i-d_j} = 1) = \frac{p(e_{l_i-d_j} = 1 | m_\delta) p(m_\delta)}{p(e_{l_i-d_j} = 1)} \quad (4)$$

$$p(m_\delta | e_{l_i-d_j} = 0) = \frac{p(e_{l_i-d_j} = 0 | m_\delta) p(m_\delta)}{p(e_{l_i-d_j} = 0)}, \quad (5)$$

where $p(e_{l_i-d_j}=1|m_\delta)$ and $p(e_{l_i-d_j}=0|m_\delta)$ represent the conditional probability of whether the lncRNA node l_i is connected to the disease node d_j or not, respectively, and m_δ is one of the common neighboring nodes between l_i and d_j in GN_1 . Thus, $p(e_{l_i-d_j}=1|m_\delta)$ and $p(e_{l_i-d_j}=0|m_\delta)$ are calculated via the following formulas:

$$p(e_{l_i-d_j} = 1 | m_\delta) = \frac{N_{m_\delta}^+}{N_{m_\delta}^+ + N_{m_\delta}^-} \quad (6)$$

$$p(e_{l_i-d_j} = 0 | m_\delta) = \frac{N_{m_\delta}^-}{N_{m_\delta}^+ + N_{m_\delta}^-}, \quad (7)$$

- 3 where $N_{m_\delta}^+$ and $N_{m_\delta}^-$ denote the number of known and unknown associations between lncRNAs and
- 4 diseases whose common neighbors include m_δ respectively.

Hence, from Equations (4) and (5), Equation (3) can be modified as follows:

$$S1(l_i, d_j) = \frac{p(e_{l_i-d_j} = 1)}{p(e_{l_i-d_j} = 0)} \prod_{m_\delta \in CN(l_i, d_j)} \frac{p(e_{l_i-d_j} = 0) p(e_{l_i-d_j} = 1 | m_\delta)}{p(e_{l_i-d_j} = 1) p(e_{l_i-d_j} = 0 | m_\delta)}. \quad (8)$$

Moreover, given any two nodes l_i and d_j in GN_1 , the value of $\frac{p(e_{l_i-d_j}=1)}{p(e_{l_i-d_j}=0)}$ is a constant, which we denote as ϕ_m for convenience. Additionally, for each common neighboring node between l_i and d_j in GN_1 , let N_l denote the number of lncRNAs directly related to m_δ , and N_d denote the number of diseases directly related to m_δ . Then, $N_{m_\delta}^+ + N_{m_\delta}^- = N_l \times N_d$, and hence, Equation (3) can further be modified as follows:

$$S1(l_i, d_j) = \phi_m \prod_{m_\delta \in CN(l_i, d_j)} \phi_m^{-1} \frac{N_{m_\delta}^+}{N_{m_\delta}^-}. \quad (9)$$

Considering that $N_{m_\delta}^+$ may equal zero, we will introduce the Laplace calibration to guarantee that the value of $S1(l_i, d_j)$ will not be zero:

$$S1(l_i, d_j) = \phi_m \prod_{m_\delta \in CN(l_i, d_j)} \phi_m^{-1} \frac{N_{m_\delta}^+ + 1}{N_{m_\delta}^- + 1}. \quad (10)$$

Furthermore, by introducing the logarithmic function for standardization, for any given lncRNA node l_i and disease node d_j in GN_1 , we can finally define the probability of a potential association existing between them as:

$$S1'(l_i, d_j) = \frac{\log(S1(l_i, d_j))}{\lambda}, \quad (11)$$

- 5 where, λ is a constant utilized for normalization.

6 1.2. Method for Applying the Naïve Bayesian Theory to GN_2

- 7 In the same manner as described in section 1.1, for any given lncRNA node l_i and
- 8 disease node d_j in GN_2 , we construct the set consisting of all common neighboring nodes,
- 9 $CN'(l_i, d_j) = \{m_1, m_2, \dots, m_h, g_1, g_2, \dots, g_u\}$. Then, the posterior probabilities of $p'(e_{l_i-d_j}=1|CN'(l_i, d_j))$

10 and $p'(e_{l_i-d_j}=0|CN'(l_i, d_j))$, representing whether the node l_i is connected to d_j in GN_2 or not,
 11 respectively, can be described as follows:

$$\begin{aligned} p'(e_{l_i-d_j}=1|CN'(l_i, d_j)) &= \frac{p'(m_1, m_2, \dots, m_h, g_1, g_2, \dots, g_u | e_{l_i-d_j}=1) p'(e_{l_i-d_j}=1)}{p'(CN'(l_i, d_j))} \\ &= \frac{p'(e_{l_i-d_j}=1)}{p'(CN'(l_i, d_j))} \prod_{m_\alpha \in CN'(l_i, d_j)} p'(m_\alpha | e_{l_i-d_j}=1) \times \\ &\quad \prod_{g_\beta \in CN'(l_i, d_j)} p'(g_\beta | e_{l_i-d_j}=1) \times \prod_{m_{\bar{\alpha}}, g_{\bar{\beta}} \in CN'(l_i, d_j)} p'(m_{\bar{\alpha}}, g_{\bar{\beta}} | e_{l_i-d_j}=1) \end{aligned} \quad (12)$$

$$\begin{aligned} p'(e_{l_i-d_j}=0|CN'(l_i, d_j)) &= \frac{p'(m_1, m_2, \dots, m_h, g_1, g_2, \dots, g_u | e_{l_i-d_j}=0) p'(e_{l_i-d_j}=0)}{p'(CN'(l_i, d_j))} \\ &= \frac{p'(e_{l_i-d_j}=0)}{p'(CN'(l_i, d_j))} \prod_{m_\alpha \in CN'(l_i, d_j)} p'(m_\alpha | e_{l_i-d_j}=0) \times \\ &\quad \prod_{g_\beta \in CN'(l_i, d_j)} p'(g_\beta | e_{l_i-d_j}=0) \times \prod_{m_{\bar{\alpha}}, g_{\bar{\beta}} \in CN'(l_i, d_j)} p'(m_{\bar{\alpha}}, g_{\bar{\beta}} | e_{l_i-d_j}=0), \end{aligned} \quad (13)$$

12 where $p'(m_\alpha | e_{l_i-d_j}=1)$ and $p'(m_\alpha | e_{l_i-d_j}=0)$ are the conditional probabilities of whether a node
 13 m_α belongs to $CN'(l_i, d_j)$; they represent the possibilities of whether the node m_α is a common
 14 neighboring node between l_i and d_j in GN_2 or not, respectively. $p'(g_\beta | e_{l_i-d_j}=1)$ and $p'(g_\beta | e_{l_i-d_j}=0)$
 15 are the conditional probabilities of whether a node g_β belonging to $CN'(l_i, d_j)$; they represent the
 16 possibilities of whether the node g_β is a common neighboring node between l_i and d_j in GN_2 or not,
 17 respectively.

Following the example of Equations (4) and (5), we have:

$$p'(m_\alpha | e_{l_i-d_j}=1) = \frac{p'(e_{l_i-d_j}=1 | g_\alpha) p'(g_\alpha)}{p'(e_{l_i-d_j}=1)} \quad (14)$$

$$p'(g_\alpha | e_{l_i-d_j}=0) = \frac{p'(e_{l_i-d_j}=0 | g_\alpha) p'(g_\alpha)}{p'(e_{l_i-d_j}=0)} \quad (15)$$

$$p'(g_\beta | e_{l_i-d_j}=1) = \frac{p'(e_{l_i-d_j}=1 | g_\beta) p'(g_\beta)}{p'(e_{l_i-d_j}=1)} \quad (16)$$

$$p'(g_\beta | e_{l_i-d_j}=0) = \frac{p'(e_{l_i-d_j}=0 | g_\beta) p'(g_\beta)}{p'(e_{l_i-d_j}=0)} \quad (17)$$

$$p'(m_{\bar{\alpha}}, g_{\bar{\beta}} | e_{l_i-d_j}=1) = \frac{p'(e_{l_i-d_j}=1 | m_{\bar{\alpha}}, g_{\bar{\beta}}) p'(m_{\bar{\alpha}}, g_{\bar{\beta}})}{p'(e_{l_i-d_j}=1)} \quad (18)$$

$$p'(m_{\bar{\alpha}}, g_{\bar{\beta}} | e_{l_i-d_j}=0) = \frac{p'(e_{l_i-d_j}=0 | m_{\bar{\alpha}}, g_{\bar{\beta}}) p'(m_{\bar{\alpha}}, g_{\bar{\beta}})}{p'(e_{l_i-d_j}=0)}, \quad (19)$$

where $p'(e_{l_i-d_j}=1 | m_{\bar{\alpha}}, g_{\bar{\beta}})$ and $p'(e_{l_i-d_j}=0 | m_{\bar{\alpha}}, g_{\bar{\beta}})$ represent the conditional probability of whether the lncRNA node l_i is connected to the disease node d_j or not respectively, while both $m_{\bar{\alpha}}$ and $g_{\bar{\beta}}$ are common neighboring nodes between l_i and d_j in GN_2 . Moreover, let $N_{m_{\bar{\alpha}}, g_{\bar{\beta}}}^+$ and $N_{m_{\bar{\alpha}}, g_{\bar{\beta}}}^-$ denote the number of known and unknown associations between l_i and d_j in GN_2 , respectively, conditional on $m_{\bar{\alpha}}$ and $g_{\bar{\beta}}$ being common neighboring nodes between l_i and d_j in GN_2 and $m_{\bar{\alpha}}-g_{\bar{\beta}}$ is an miRNA-gene

pair. In addition, for any given miRNA–gene pair $m_{\bar{\alpha}}-g_{\bar{\beta}}$, let N'_{l_1} denote the number of lncRNAs directly related to $m_{\bar{\alpha}}-g_{\bar{\beta}}$, and N'_{d_1} denote the number of diseases directly related to $m_{\bar{\alpha}}-g_{\bar{\beta}}$; then, $N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^+ + N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^- = N'_{l_1} \times N'_{d_1}$ and $p'(e_{l_i-d_j}=1|m_{\bar{\alpha}},g_{\bar{\beta}}) + p'(e_{l_i-d_j}=0|m_{\bar{\alpha}},g_{\bar{\beta}})=1$. Therefore, we have:

$$p'(e_{l_i-d_j}=1|m_{\bar{\alpha}},g_{\bar{\beta}}) = \frac{N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^+}{N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^+ + N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^-}. \quad (20)$$

$$p'(e_{l_i-d_j}=0|m_{\bar{\alpha}},g_{\bar{\beta}}) = \frac{N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^-}{N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^+ + N_{m_{\bar{\alpha}},g_{\bar{\beta}}}^-}. \quad (21)$$

As illustrated in Section 1.1, we can define the probability of a potential association existing between l_i and d_j in GN_2 as follows:

$$S2(l_i, d_j) = \phi_m \prod_{m_{\alpha} \in CN'(l_i, d_j)} \prod_{g_{\beta} \in CN'(l_i, d_j)} \prod_{m_{\bar{\alpha}}, g_{\bar{\beta}} \in CN'(l_i, d_j)} \phi_m^{-3} \frac{(N_{m_{\alpha}}^+ + 1)(N_{g_{\beta}}^+ + 1)(N_{m_{\bar{\alpha}}, g_{\bar{\beta}}}^+ + 1)}{(N_{m_{\alpha}}^- + 1)(N_{g_{\beta}}^- + 1)(N_{m_{\bar{\alpha}}, g_{\bar{\beta}}}^- + 1)}, \quad (22)$$

where $N_{m_{\alpha}}^+$ and $N_{m_{\alpha}}^-$ denote the number of known and unknown associations between l_i and d_j in GN_2 respectively, conditional on m_{α} being a common neighboring node between l_i and d_j . In addition, $N_{g_{\beta}}^+$ and $N_{g_{\beta}}^-$ represent the number of known and unknown associations between l_i and d_j in GN_2 respectively, conditional on g_{β} being a common neighboring node between l_i and d_j . From the above descriptions, we find $N_{m_{\alpha}}^+ + N_{m_{\alpha}}^- = N'_{l_2} \times N'_{d_2}$ and $N_{g_{\beta}}^+ + N_{g_{\beta}}^- = N'_{l_3} \times N'_{d_3}$, where N'_{l_2} denotes the number of lncRNAs directly related to the node m_{α} in GN_2 , N'_{d_2} denotes the number of diseases directly related to the node m_{α} in GN_2 , N'_{l_3} denotes the number of lncRNAs directly related to the node g_{β} in GN_2 , and N'_{d_3} denotes the number of diseases directly related to the node in GN_2 .

Finally, following the example of Equation (11), we can finally define the probability of a potential association existing between l_i and d_j in GN_2 as follows:

$$S2'(l_i, d_j) = \frac{\log(S2(l_i, d_j))}{\lambda}. \quad (23)$$