

Predicting Compound-Protein Interactions by Multi-Task Learning

Songpeng Zu

January 28, 2015

It is a key problem to evaluate the compound-protein interactions (CPIs) in the early stage of drug discovery and drug design. Many *in silico* methods have been used to predict compound-protein interactions due to the high cost and high risk experiments.

In this paper, we develop a Bayesian Hierarchical model to predict compound-protein interactions. The main innovations are : a) we integrate both chemoinformation and genomic information, i.e., large-scale compound-protein interactions information, compounds 2D structure fingerprints, chemo-physics information and several levels protein information, such as proteins family information, proteins sequence information, functional information and protein-protein interactions network information. b) instead of predicting whether or not they interact like most papers do, here we can predict the CPIs potencies. c) last but most importantly, we predict compound-protein interactions by ligand-based approach, but not independently. That means we incorporate other proteins chemoinformation into the current learning procedure, namely which kind of compounds can interact with a specific proteins, by considering proteins similarities from different levels.

CONTENTS

1	Introduction	3
2	Method and Materials	4
2.1	Materials	4
2.2	Method	4
3	Result	6
3.1	The performance of different chemical features	6
3.2	Compare with modeling on a single protein	6
3.3	Compare with modeling on a protein family	6
3.4	Chemical substructures against different proteins	6
3.5	New scaffold discovery against several proteins from traditional Chinese Medicine herbs	6
3.6	Discovery the compounds targeting RIP3 or PSH	6
4	Discussion	6
5	Append	7

1 INTRODUCTION

- New chemical scaffold

2 METHOD AND MATERIALS

2.1 MATERIALS

2.2 METHOD

Suppose we have m groups of data. In each group, we have the data $\mathcal{D}_j = \{\mathbf{X}_j, \mathbf{y}_j\}$, $j = 1, \dots, m$, and $\mathbf{X}_j \in \mathbb{R}^{d \times n_j}$. Then we have

$$\mathbf{y}_j \sim \mathcal{N}(\mathbf{X}_j^T \omega_j, \sigma_y^2 \mathbf{I}) \quad (2.1)$$

Since different groups data may share similar or common structures, we assume the parameters of ω_j have the same mean on the prior distribution.

$$\omega_j \sim \mathcal{N}(\omega_*, \sigma_j^2 \mathbf{I}) \quad (2.2)$$

In which,

$$\omega_* \sim \mathcal{N}(\mu, \sigma_*^2 \mathbf{I}) \quad (2.3)$$

Suppose, for simplicity, that $\mu = \mathbf{0}$, $p(\sigma_y^2) \propto 1$, and that σ_j^2 and σ_* are fixed. Let $\Theta = \{\omega_j, j = 1, \dots, m, \omega_*, \sigma_y^2\}$. We have

$$\begin{aligned} \mathcal{L}_{hier}(\mathcal{D}; \Theta) &= \mathcal{L}_{orig}(\mathcal{D} | \Theta) + \log p(\Theta) \\ &= \sum_j \left(\log p(\mathcal{D}_j | \omega_j) - \frac{\|\omega_j - \omega_*\|^2}{2\sigma_j^2} \right) - \frac{\|\omega_*\|^2}{2\sigma_*^2} - \overbrace{\sum_j \frac{d}{2} \log(2\pi\sigma_j^2) - \frac{d}{2} \log(2\pi\sigma_*^2)}^{Const} \\ &= \sum_j \left(-\frac{\|\mathbf{y}_j - \mathbf{X}_j^T \omega_j\|^2}{2\sigma_y^2} - \frac{\|\omega_j - \omega_*\|^2}{2\sigma_j^2} \right) - \frac{\|\omega_*\|^2}{2\sigma_*^2} - \sum_j \frac{n_j}{2} \log(2\pi\sigma_y^2) \\ &\quad - \overbrace{\sum_j \frac{d}{2} \log(2\pi\sigma_j^2) - \frac{d}{2} \log(2\pi\sigma_*^2)}^{Const} \end{aligned} \quad (2.4)$$

Then we can get

$$\begin{aligned} \frac{\partial \mathcal{L}_{hier}(\mathcal{D}; \Theta)}{\partial \omega_j} &= -\frac{1}{2\sigma_y^2} \frac{\|\mathbf{y}_j - \mathbf{X}_j^T \omega_j\|^2}{\partial \omega_j} - \frac{1}{2\sigma_j^2} \frac{\|\omega_j - \omega_*\|^2}{\partial \omega_j} \\ &= \frac{\mathbf{X}_j \mathbf{y}_j}{\sigma_y^2} + \frac{\omega_*}{\sigma_j^2} - \left(\frac{\mathbf{X}_j \mathbf{X}_j^T}{\sigma_y^2} + \frac{1}{\sigma_j^2} \mathbf{I} \right) \omega_j \end{aligned} \quad (2.5)$$

$$\frac{\partial \mathcal{L}_{hier}(\mathcal{D}; \Theta)}{\partial \omega_*} = -\sum_j \frac{\omega_* - \omega_j}{\sigma_j^2} - \frac{\omega_*}{\sigma_*^2} \quad (2.6)$$

$$\frac{\partial \mathcal{L}_{hier}(\mathcal{D}; \Theta)}{\partial \sigma_y^2} = \frac{\sum_j \|\mathbf{y}_j - \mathbf{X}_j^T \omega_j\|^2}{2(\sigma_y^2)^2} - \frac{n}{2\sigma_y^2} \quad (2.7)$$

where n is the total number of samples in all the groups.

In this way, L-BFGS-B method is used to get the MAP estimations of these parameters following the function of $\mathcal{L}_{hier}(\mathcal{D}; \omega)$ and their derivations ¹.

Also we can derived the conditional posterior distribution of these parameters.

$$p(\omega_j | \omega_{[-j]}, \sigma_y^2) \sim \mathcal{N}(\mu_j, \Sigma_j) \quad (2.8)$$

$$\mu_j = \Sigma_j \eta_j, \quad \eta_j = \frac{\mathbf{X}_j \mathbf{y}_j}{\sigma_y^2} + \frac{\omega_*}{\sigma_j^2}, \quad \Sigma_j^{-1} = \frac{\mathbf{I}}{\sigma_j^2} + \frac{\mathbf{X}_j^T \mathbf{X}_j}{\sigma_y^2}$$

$$p(\omega_* | \omega_{[-*]}, \sigma_y^2) \sim \mathcal{N}(\mu_*, \Sigma_*) \quad (2.9)$$

$$\mu_* = \Sigma_* \eta_*, \quad \eta_* = \sum_j \frac{\omega_j}{\sigma_j^2}, \quad \Sigma_*^{-1} = \sum_j \frac{\mathbf{I}}{\sigma_j^2} + \frac{\mathbf{I}}{\sigma_*^2}$$

$$p(\sigma_y^2 | \omega) \sim Inv - \chi^2(v, s^2) \quad (2.10)$$

$$v = n - 2, \quad s^2 = \frac{\sum_j \|\mathbf{y}_j - \mathbf{X}_j^T \omega_j\|^2}{n - 2}$$

Then we can also use Gibbs Samplers method to get the final results.

If m equals to 1, our model is reduced into the ridge regression model.

$$\omega_j \sim \mathcal{N}(\mathbf{0}, (\sigma_j^2 + \sigma_*^2) \mathbf{I}) \quad (2.11)$$

This can be derived based on the relationship

$$E(x) = E(E(x|y)), \quad Var(x) = E(Var(x|y)) + Var(E(x|y)) \quad (2.12)$$

and the marginal distribution is normal distribution if the hyperprior of mean and the conditional distribution is normal.

¹Note that $\sigma_y^2 > 0$, this is an bound-constrained optimization. L-BFGS-B is memory-limited ,especially used in bound-constrained, and suitable for quite large-dimensional optimization.

3 RESULT

3.1 THE PERFORMANCE OF DIFFERENT CHEMICAL FEATURES

3.2 COMPARE WITH MODELING ON A SINGLE PROTEIN

3.3 COMPARE WITH MODELING ON A PROTEIN FAMILY

3.4 CHEMICAL SUBSTRUCTURES AGAINST DIFFERENT PROTEINS

3.5 NEW SCAFFOLD DISCOVERY AGAINST SEVERAL PROTEINS FROM TRADITIONAL CHINESE MEDICINE HERBS

3.6 DISCOVERY THE COMPOUNDS TARGETING RIP3 OR PSH

4 DISCUSSION

REFERENCE

5 APPEND

Suppose there are M different groups of data, i.e.,

$$\left\{ (y_i^{(g)}, \mathbf{x}_i^{(g)}) \right\}_{i=1}^{n^{(g)}}; 1 \leq g \leq M$$

In our work, we have M proteins, and for each protein g ($1 \leq g \leq M$), there are n_g compounds target it. $y_i^{(g)}$ represents the potency of the compound i against the protein g . $\mathbf{x}_i^{(g)}$, a binary vector which has D dimensions, represents the compound i 's two-dimensional chemical structure fingerprints in the protein g . Here we use a linear model to describe the relationship between the potency and chemical structure fingerprints³. All the vectors, if no special explanations, are column vectors.

$$y_i^{(g)} = \mu^{(g)} + \mathbf{x}_i^{(g)T} \cdot \boldsymbol{\beta}^{(g)} + \varepsilon \quad (5.1)$$

In which, $\mu^{(g)}$ is the average effect⁴, $\boldsymbol{\beta}^{(g)}$ is the coefficients and ε is the systematic random error.

$$\varepsilon \sim N(0, \tau^{-1}) \quad (5.2)$$

$$\tau \sim \text{Gamma}(\kappa_1, \kappa_2) \quad (5.3)$$

And

$$\mu^{(g)} | \tau \sim N(0, \sigma_\mu^2 \tau^{-1}) \quad (5.4)$$

Considering the fact that proteins having similar pharmacological properties tend to interact with structure similar compounds, we construct a hierarchical model to fit this phenomenon.

Let

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_{11} & \beta_{12} & \dots & \beta_{1D} \\ \beta_{21} & \beta_{22} & \dots & \beta_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{M1} & \beta_{M2} & \dots & \beta_{MD} \end{pmatrix}$$

Then $\boldsymbol{\beta}^{(g)}$ corresponds to the g th row of $\boldsymbol{\beta}$ matrix and let $\boldsymbol{\beta}_{(d)}$ corresponds to the d th column of $\boldsymbol{\beta}$ matrix. We assume that

$$\boldsymbol{\beta}_{(d)} \sim \text{MVN}(\mathbf{0}, \sigma^2 \tau^{-1} \Sigma_\beta) \quad (5.5)$$

In which, σ^2 is a scale parameter, and follows a non-informative prior,

$$\text{Pr}(\sigma^2) \propto \frac{1}{\sigma^2} \quad (5.6)$$

Σ_β can be treated as a kind of covariance matrix. In our model,

$$\Sigma_\beta = \sum_{1 \leq n \leq N} \omega_n \Omega^{(n)}. \quad (5.7)$$

²The variance of potencies might be too large, should take log.

³Further complex function can be involved, like Gauss Process

⁴If y is centered, we can ignore the μ .

In which, Σ^n is a similarity Matrix between different proteins. Since we have several ways to define their similarities, such as sequence similarity, GO functional similarity, the distances in the protein-protein interaction networks and so on. So n corresponds to the n th definition of the protein similarity and in total we have N ways. ω_n is the weight for Σ^n .⁵

Then we can derive the posterior distribution as followed, let $\Theta = \{\tau, \sigma^2, \beta, \mu\}$. Then,

$$\begin{aligned}
Pr(\Theta|\mathbf{y}, \mathbf{X}) &\propto Pr(\tau)Pr(\sigma^2)Pr(\beta|\tau, \sigma^2)Pr(\mu|\tau)Pr(\mathbf{y}|\mathbf{X}, \Theta) \\
&\propto \tau^{\kappa_1-1} e^{-\tau\kappa_2} \frac{1}{\sigma^2} \prod_{d=1}^D \frac{1}{\sqrt{(2\pi)^M |\sigma^2 \Sigma_\beta \tau^{-1}|}} \exp \left\{ -\frac{1}{2} \beta_{(d)}^T (\sigma^2 \tau^{-1} \Sigma_\beta)^{-1} \beta_{(d)} \right\} \\
&\cdot \prod_{g=1}^M \frac{1}{\sqrt{2\pi\sigma_\mu^2 \tau^{-1}}} \exp \left\{ -\frac{1}{2} \frac{\mu^{(g)^2}}{\sigma_\mu^2 \tau^{-1}} \right\} \cdot \prod_{g=1}^M \prod_{i=1}^{n^{(g)}} \frac{1}{\sqrt{2\pi\tau^{-1}}} \exp \left\{ -\frac{1}{2} \frac{(y_i^{(g)} - \mu^{(g)} - \mathbf{x}_i^{(g)T} \beta^{(g)})^2}{\tau^{-1}} \right\}
\end{aligned} \tag{5.8}$$

Then, we can derive the conditional posterior distribution for these parameters.

The conditional posterior distribution of $\beta_{(d)}$ follows multivariable normal distribution,

$$\begin{aligned}
Pr(\beta_{(d)}|\mathbf{y}, \mathbf{X}, \beta_{[-d]}, \tau, \sigma^2, \mu) &\propto Pr(\beta_{(d)}|\sigma^2, \tau)Pr(\mathbf{y}|\mathbf{X}, \Theta) \\
&\propto \exp \left\{ -\frac{1}{2} \beta_{(d)}^T (\sigma^2 \tau^{-1} \Sigma_\beta^{-1}) \beta_{(d)} \right\} \\
&\cdot \exp \left\{ -\frac{\tau}{2} \beta_d^T \Psi_d \beta_d + \tau \boldsymbol{\varphi}^T \beta_{(d)} + Const \right\}
\end{aligned} \tag{5.9}$$

The conditional posterior distribution of τ follows Gamma distribution,

$$\begin{aligned}
Pr(\tau|\mathbf{y}, \mathbf{X}, \beta, \sigma^2, \mu) &\propto Pr(\tau)Pr(\beta|\tau, \sigma^2)Pr((\mu)|\tau)Pr(\mathbf{y}|\mathbf{X}, \Theta) \\
&\propto \tau^{\kappa_1 + \frac{D \cdot M}{2} + \frac{M}{2} + \mathbf{M} \cdot \mathbf{n}^{(g)}/2} \exp \left\{ -\tau\kappa_2 - \tau \sum_{d=1}^D \frac{1}{2} \beta_{(d)}^T (\sigma^2 \Sigma_\beta)^{-1} \beta_{(d)} - \tau \sum_{g=1}^M \frac{1}{2} \frac{\mu^{(g)^2}}{\sigma_\mu^2} \right. \\
&\quad \left. - \tau \sum_{g=1}^M \sum_{i=1}^{n^{(g)}} \frac{1}{2} (y_i^{(g)} - \mu^{(g)} - \mathbf{x}_i^{(g)T} \beta^{(g)})^2 \right\}
\end{aligned} \tag{5.10}$$

⁵We can also maximum similarity or the correlation among them.