

Supplemental file

Network as a biomarker: A novel network-based sparse Bayesian machine for pathway-driven drug response prediction

Qi Liu^{1,2}, Louis J. Muglia^{2,3} and Lei Frank Huang^{1,2,4*}

¹ Brain Tumor Center, Division of Experimental Hematology and Cancer Biology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH 45229, USA

² Department of Pediatrics, College of Medicine, University of Cincinnati, Cincinnati, OH 45229, USA

³ Division of Human Genetics, Cincinnati Children's Hospital Medical Center, Cincinnati, OH 45229, USA

⁴ Department of Information Science, School of Mathematical Sciences and LAMA, Peking University, Beijing 100871, China

* Correspondence: Frank.Huang@cchmc.org; Tel.: +15135171084

1 Approximate Bayesian Inference for parameter estimation in NBSBM using Expectation Propagation (EP) algorithm

Considering the labeling errors ε , given $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, $\mathbf{y} = (y_1, \dots, y_n)$ and ε , the likelihood can be written as (1)

$$p(\mathbf{y}|\beta, \varepsilon, \mathbf{X}) = \prod_{i=1}^n p(y_i|\beta, \varepsilon, \mathbf{x}_i) = \prod_{i=1}^n [\varepsilon(1 - \Phi(y_i\beta^T \mathbf{x}_i)) + (1 - \varepsilon)\Phi(y_i\beta^T \mathbf{x}_i)] = \prod_{i=1}^n [\varepsilon + (1 - 2\varepsilon)\Phi(y_i\beta^T \mathbf{x}_i)] \quad (1)$$

Where Φ is the Heaviside step function and it is defined by equation (2)

$$\Phi(y_i\beta^T \mathbf{x}_i) = \lim_{k \rightarrow \infty} \frac{1}{1 + e^{-2k(y_i\beta^T \mathbf{x}_i)}} \quad (2)$$

If we only consider the sparse solution for β . Herein we introduce a new binary hidden variable $\mathbf{z} = \{z_0, z_1, z_2, \dots, z_d\} \in \{0, 1\}^d$. z_i takes 0 if the i^{th} component of β_{true} is 0 and z_i takes 1 otherwise. Assuming \mathbf{z} is given, the probability density of β is shown in equation (3)

$$p(\beta|\mathbf{z}) = \prod_{i=1}^d p(\beta_i|z_i) = \prod_{i=0}^d [\mathcal{N}(\beta_i, 0, \sigma_i^2)^{z_i} (\delta(\beta_i))^{(1-z_i)}] \quad (3)$$

where $p(\beta_i|z_i)$ is a Spike and Slab prior. $\mathcal{N}(\beta_i, 0, \sigma_i^2)$ represents Gaussian density function with 0 mean and σ_i^2 variance, $\delta(\beta_i)$ is an impulse function which has a probability of 1 on β_i and 0 elsewhere. To complete the specification of the prior for β at zero, we assume that a network that encodes the dependencies between the gene features are known. Given a specific cancer signaling network $G = (V, E)$ whose vertices $V = \{0, 1, \dots, d\}$ correspond to the proteins and whose edges, E Equation (4) shows the prior probability for \mathbf{z} given G which is given by a Markov random field (MRF) model

$$p(\mathbf{z}|G, \lambda, \gamma) = \frac{1}{Z} \exp\left(cz_0 + \lambda \sum_{i=1}^d z_i + \gamma \sum_{\{u,v\} \in E} \left(\frac{z_u}{\sqrt{d_u}} - \frac{z_v}{\sqrt{d_v}}\right)^2 w(u, v)\right) = \frac{1}{Z} \exp(cz_0 + \lambda \sum_{i=1}^d z_i) \exp\left(\gamma \sum_{\{u,v\} \in E} \left(\frac{z_u}{\sqrt{d_u}} - \frac{z_v}{\sqrt{d_v}}\right)^2 w(u, v)\right) \quad (4)$$

In equation (2), Z is a normalization constant and $\lambda \in \mathbb{R}$ controls the sparsity. $\gamma \geq 0$ determines the sum of square difference between z_u and z_v that are linked in the input network G , $w(u, v)$ is the weight between proteins z_u and z_v . In fact, if we assume,

$$L(u, v) = \begin{cases} 1 - \frac{w(u,v)}{d_u}, & \text{if } u = v \text{ and } d_u \neq 0, \\ \frac{-w(u,v)}{\sqrt{d_u d_v}}, & \text{if } u \text{ and } v \text{ are adjacent,} \\ 0, & \text{othersize.} \end{cases} \quad (5)$$

then

$$p(\mathbf{z}|G, \lambda, \gamma) = \frac{1}{Z} \exp(cz_0 + \lambda|\mathbf{z}|) \exp(\gamma \mathbf{z}^T L \mathbf{z}) \quad (6)$$

Furthermore, we assume the prior of ε as

$$p(\varepsilon) = \text{Beta}(\varepsilon, a_0, b_0) = \frac{1}{B(a_0, b_0)} \varepsilon^{a_0-1} (1 - \varepsilon)^{b_0-1} \quad (7)$$

where $B(a_0, b_0)$ represents beta function with parameters a_0 and b_0 . Under the assumption above, we can use Bayesian theorem to compute the posterior distribution of the model parameters β and ε given the training data \mathbf{X} and \mathbf{y} . Given the specific cancer signaling network G and the model hyper-parameters λ and γ , the posterior is given by

$$p(\beta, \varepsilon | \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) = \frac{\sum_{\mathbf{z}} p(\mathbf{y} | \beta, \varepsilon, \mathbf{X}) p(\beta | \mathbf{z}) p(\mathbf{z} | G, \lambda, \gamma) p(\varepsilon)}{p(\mathbf{y} | \mathbf{X}, G, \lambda, \gamma)} \quad (8)$$

If given a new unclassified sample x^{test} , we can determine its classification labels y^{test} by probability as shown in equation (9)

$$p(y^{\text{test}} | X^{\text{test}}, y, \mathbf{X}, G, \lambda, \gamma) = \iint p(y^{\text{test}} | \beta, \varepsilon, x^{\text{test}}) p(\beta, \varepsilon | \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) d\beta d\varepsilon \quad (9)$$

Then the relevance of the features can be quantified by the posterior of \mathbf{z} ,

$$p(\mathbf{z} | \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) = \frac{\sum_{\beta} \sum_{\varepsilon} p(\mathbf{y} | \beta, \varepsilon, \mathbf{X}) p(\beta | \mathbf{z}) p(\mathbf{z} | G, \lambda, \gamma) p(\varepsilon)}{p(\mathbf{y} | \mathbf{X}, G, \lambda, \gamma)} \quad (10)$$

In specific, the relevance of the i -th feature to the classification result is a value between 0 and 1 and is given by the marginal probability $p(\mathbf{z} | \mathbf{y}, \mathbf{X}, G, \lambda, \gamma)$ with $\mathbf{z} = \mathbf{1}$. The higher the value, the more relevant of this gene with respect to the classification result. The joint probability distributions of model parameters and hidden variables are given as follows:

$$p(\beta, \varepsilon, \mathbf{z} | \mathbf{X}, G, \lambda, \gamma) = p(\mathbf{y} | \beta, \varepsilon, \mathbf{X}) p(\beta | \mathbf{z}) p(\mathbf{z} | G, \lambda, \gamma) p(\varepsilon) \quad (11)$$

It can be written as the product of $N + |E| + 3$ probabilities in equation (11) according to the assumption of independence.

$$p(\beta, \varepsilon, \mathbf{z}, \mathbf{y}|\mathbf{X}, G, \lambda, \gamma) = [\prod_{i=1}^n p(y_i|\beta, \varepsilon, \mathbf{x}_i)] [\prod_{i=0}^d p(\beta_i|z_i)] p(\mathbf{z}|G, \lambda, \gamma) p(\varepsilon) = \prod_{i=1}^{n+|E|+3} t_i(\beta, \varepsilon, \mathbf{z}) = q(\beta, \varepsilon, \mathbf{z}) \quad (12)$$

Where $|E|$ refers to the number of edges in graph G . The first n terms of $t_i(\beta, \varepsilon, \mathbf{z})$ denote the likelihood $p(y_i|\beta, \varepsilon, \mathbf{x}_i)$, while $t_{n+1}(\beta, \varepsilon, \mathbf{z})$, $\prod_{i=n+2}^{n+|E|+2} t_i(\beta, \varepsilon, \mathbf{z})$ and $t_{n+|E|+3}(\beta, \varepsilon, \mathbf{z})$ represent $p(\beta|\mathbf{z})$, $p(\mathbf{z}|G, \lambda, \gamma)$ and $p(\varepsilon)$ respectively. According to the expectation propagation algorithm, we use \tilde{t}_i as the estimation of t_i and get (13)

$$\prod_{i=1}^{n+|E|+3} t_i(\beta, \varepsilon, \mathbf{z}) \approx \prod_{i=1}^{n+|E|+3} \tilde{t}_i(\beta, \varepsilon, \mathbf{z}) = Q(\beta, \varepsilon, \mathbf{z}) \quad (13)$$

It is restricted that all \tilde{t}_i belong to the same exponential family of distributions, and $Q(\beta, \varepsilon, \mathbf{z})$ have the same expression with $\tilde{t}_i(\beta, \varepsilon, \mathbf{z})$ because the product of functions belonging to the same exponential family of distributions is a closure. Assume that the density function of Q after normalization is Q , which is also the approximation of the posterior distribution $p(\beta, \varepsilon, \mathbf{z}, \mathbf{y}|\mathbf{X}, G, \lambda, \gamma)$, and use $Q^{\setminus i}(\beta, \varepsilon, \mathbf{z})$ to denote the approximation of $Q(\beta, \varepsilon, \mathbf{z})$ without the term t_i as shown in (8)

$$Q^{\setminus i}(\beta, \varepsilon, \mathbf{z}) = \prod_{j \neq i} \tilde{t}_j(\beta, \varepsilon, \mathbf{z}) = \frac{Q(\beta, \varepsilon, \mathbf{z})}{\tilde{t}_i(\beta, \varepsilon, \mathbf{z})} \quad (14)$$

a general workflow of the expectation propagation algorithm for the sparse Bayesian classifier can be given as follows,

1. Initialize all \tilde{t}_i and posterior distribution Q ;
2. Repeat the following steps until all \tilde{t}_i converge.

(a) Select one \tilde{t}_i that needs to be changed and calculate $Q^{\setminus i}$: $Q^{\setminus i} = Q / \tilde{t}_i$

(b) Update the value of Q to minimize the Kullback–Leibler(KL) divergence between $t_i Q^{\setminus i}$ and $\tilde{t}_i Q^{\setminus i}$.

(c) Recalculate $\tilde{t}_i = Q^{\text{new}} / Q^{\setminus i}$.

3. Estimate model parameters.

In fact, according to (1), (3), (4) and (5), we can approximate \tilde{t}_i based on function (14)

$$\tilde{t}_i(\beta, \varepsilon, \mathbf{z}) = \tilde{s}_i \varepsilon^{\tilde{a}_i} (1 - \varepsilon)^{\tilde{b}_i} \prod_{j=0}^d \exp\left(-\frac{1}{2\tilde{v}_{ij}} (\beta_j - \tilde{m}_{ij})^2\right) (z_i \tilde{c}_{ij} + (1 - z_i) \tilde{d}_{ij}) \quad (15)$$

Where $\tilde{\mathbf{m}}_i = (\tilde{m}_{i0}, \dots, \tilde{m}_{id})^T$, $\tilde{\mathbf{v}}_i = (\tilde{v}_{i0}, \dots, \tilde{v}_{id})^T$, $\tilde{\mathbf{c}}_i = (\tilde{c}_{i0}, \dots, \tilde{c}_{id})^T$, $\tilde{\mathbf{d}}_i = (\tilde{d}_{i0}, \dots, \tilde{d}_{id})^T$ and $\tilde{c}_i = 1 - \tilde{d}_i$. $\tilde{\mathbf{a}}_i$ and $\tilde{\mathbf{b}}_i$ are free parameters and \tilde{s}_i is a constant to ensure that $\tilde{t}_i Q^{\setminus i}$ and $t_i Q^{\setminus i}$ get the same value when integrating. According to the previous assumption that all \tilde{t}_i belong to the same exponential family of distributions, Q and \tilde{t}_i have the same form and we can assume that Q can be expressed as shown in (15).

$$Q(\beta, \varepsilon, \mathbf{z}) = \text{Beta}(\varepsilon|a, b) \prod_{j=0}^d \mathcal{N}(\beta_j|m_j, v_j) \text{Bern}(z_j|\rho_i) \quad (16)$$

Formula (15) has the same form as (16), and

$$\text{Bern}(z_i|\rho_i) = z_i\rho_i + (1 - z_i)(1 - \rho_i) \quad (17)$$

Where $z_i \in \{0,1\}$ and ρ_i is the probability of $z_i=1$. $\mathbf{m} = (m_0, \dots, m_d)^T$, $\mathbf{v} = (v_0, \dots, v)^T$, $\boldsymbol{\rho} = (\rho_0, \dots, \rho_d)^T$. Firstly, we initialize \mathcal{Q} and $\tilde{\tau}_i$ by setting $a = b = 1$, $m_i = 0$, $v_i = +\infty$, $\rho_i = 0.5$, $\tilde{m}_{ij} = 0$, $\tilde{v}_{ij} = +\infty$ and $\tilde{c}_{ij} = \tilde{d}_{ij} = 1$ for i in range $[1, n+|E|+3]$ and j in range $[0,d]$. Besides, as \mathcal{Q} and $\mathcal{Q}^{\setminus i}$ has the same form without approximation term $\tilde{\tau}_i$, we can make following assumption

$$\mathcal{Q}^{\setminus i}(\boldsymbol{\beta}, \varepsilon, \mathbf{z}) = \text{Beta}(\varepsilon|a^{\setminus i}, b^{\setminus i}) \prod_{j=0}^d \mathcal{N}(\beta_j|m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j|\rho_j^{\setminus i}) \quad (18)$$

$\mathbf{m}^{\setminus i} = (m_0^{\setminus i}, \dots, m_d^{\setminus i})^T$, $\mathbf{v}^{\setminus i} = (v_0^{\setminus i}, \dots, v_d^{\setminus i})^T$, $\boldsymbol{\rho}^{\setminus i} = (\rho_0^{\setminus i}, \dots, \rho_d^{\setminus i})^T$, $a^{\setminus i}$ and $b^{\setminus i}$ can be calculated based on $\mathcal{Q}^{\setminus i} = \mathcal{Q} / \tilde{\tau}_i$ and formula (15), (16).

$$\mathbf{v}^{\setminus i} = (\mathbf{v}^{-1} - \tilde{\mathbf{v}}_i^{-1})^{-1} \quad (19)$$

$$\mathbf{m}^{\setminus i} = \mathbf{m} + \mathbf{v}^{\setminus i} \circ \tilde{\mathbf{v}}_i^{-1} \circ (\mathbf{m} - \tilde{\mathbf{m}}_i) \quad (20)$$

$$\boldsymbol{\rho}^{\setminus i} = \boldsymbol{\rho} \circ \tilde{\mathbf{c}}_i^{-1} \circ (\boldsymbol{\rho} \circ \tilde{\mathbf{c}}_i^{-1} + (1 - \boldsymbol{\rho}) \circ \tilde{\mathbf{d}}_i^{-1})^{-1} \quad (21)$$

$$a^{\setminus i} = a - \tilde{a}_i \quad (22)$$

$$b^{\setminus i} = b - \tilde{b}_i \quad (23)$$

Where \circ denotes the Hadamard production and the inverse of a vector means the inverse of each component of the vector. Meanwhile we have $\tilde{\tau}_i$ satisfying (24)-(28) which can be used to update the value of $\tilde{\tau}_i$ according to the property of the exponential family functions

$$\mathbb{E}_{\tilde{\tau}_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta}] = \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta}] \quad (24)$$

$$\mathbb{E}_{\tilde{\tau}_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta} \circ \boldsymbol{\beta}] = \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta} \circ \boldsymbol{\beta}] \quad (25)$$

$$\mathbb{E}_{\tilde{\tau}_i \mathcal{Q}^{\setminus i}}[\mathbf{z}] = \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\mathbf{z}] \quad (26)$$

$$\mathbb{E}_{\tilde{\tau}_i \mathcal{Q}^{\setminus i}}[\log(\varepsilon)] = \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\log(\varepsilon)] \quad (27)$$

$$\mathbb{E}_{\tilde{\tau}_i \mathcal{Q}^{\setminus i}}[\log(1 - \varepsilon)] = \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\log(1 - \varepsilon)] \quad (28)$$

We need to update the parameters in $\tilde{\tau}_i$ according to $p(y_i|\boldsymbol{\beta}, \varepsilon, \mathbf{x}_i)$ so that $\tilde{\tau}_i$ match the constraints in (24)-(28) while minimizing the KL-divergence between $\tilde{\tau}_i \mathcal{Q}^{\setminus i}$ and $t_i \mathcal{Q}^{\setminus i}$. We can get (29) – (30) based on (1), (18), (24), (25)

$$\mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta}] = \mathbf{m}^{\setminus i} + \mathbf{v}^{\setminus i} \nabla_{\mathbf{m}} \log Z_i \quad (29)$$

$$\mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta} \circ \boldsymbol{\beta}] - \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta}] \mathbb{E}_{t_i \mathcal{Q}^{\setminus i}}[\boldsymbol{\beta}]^T = \mathbf{v}^{\setminus i} - \mathbf{v}^{\setminus i} \mathbf{v}^{\setminus i} (\nabla_{\mathbf{m}}^T \nabla_{\mathbf{m}} - 2 \nabla_{\mathbf{v}} \log Z_i) \quad (30)$$

$$Z_i = \int (\varepsilon + (1 - 2\varepsilon)\Phi(y_i \boldsymbol{\beta}^T \mathbf{x}_i)) \text{Beta}(\varepsilon|a^{\setminus i}, b^{\setminus i}) \prod_{j=1}^d \mathcal{N}(\beta_j|m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j|\rho_j^{\setminus i}) d\boldsymbol{\beta} \quad (31)$$

By substituting (29), (30) and (31) into (24) and (25), we can get

$$\mathbf{m}^{new} = \mathbf{m}^i + \mathbf{v}^i \circ \frac{(1-2\bar{\varepsilon}^i)\mathcal{N}(\lambda_i, 0, 1)}{\bar{\varepsilon}^i + (1-\bar{\varepsilon}^i)\Phi(\lambda_i)} \frac{y_i \mathbf{x}_i}{\sqrt{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i}} \quad (32)$$

$$\mathbf{v}^{new} = \mathbf{v}^i - \mathbf{v}^i \mathbf{v}^i \left(\frac{(1-2\bar{\varepsilon})\mathcal{N}(\lambda_i, 0, 1)}{\bar{\varepsilon}^i + (1-2\bar{\varepsilon}^i)\Phi(\lambda_i)} \right)^2 \frac{(y_i \mathbf{x}_i) \circ (y_i \mathbf{x}_i)}{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i} - 2 \frac{(1-2\bar{\varepsilon}^i)\mathcal{N}(\lambda_i, 0, 1)}{\bar{\varepsilon}^i + (1-2\bar{\varepsilon}^i)\Phi(\lambda_i)} \frac{y_i (\mathbf{m}^i)^T \mathbf{x}_i}{2} \frac{\mathbf{x}_i \circ \mathbf{x}_i}{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i \sqrt{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i}} \quad (33)$$

After simplifying (33), we get (34)

$$\mathbf{v}^{new} = \mathbf{v}^i - (\mathbf{v}^i \circ \mathbf{x}_i)(\mathbf{v}^i \circ \mathbf{x}_i) = \frac{y_i \alpha_i \mathbf{x}_i^T \mathbf{m}^{new}}{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i \sqrt{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i}} \quad (34)$$

In equations (32), (33) and (34), we have the following

$$\alpha_i = \frac{(1-2\bar{\varepsilon})\mathcal{N}(\lambda_i, 0, 1)}{\bar{\varepsilon}^i + (1-2\bar{\varepsilon}^i)\Phi(\lambda_i)} \quad (35)$$

$$\lambda_i = \frac{y_i (\mathbf{m}^i)^T \mathbf{x}_i}{\sqrt{\mathbf{x}_i^T \mathbf{v}^i \mathbf{x}_i}} \quad (36)$$

$$\bar{\varepsilon}^i = \frac{a^i}{a^i + b^i} \quad (37)$$

$$Z_i = \bar{\varepsilon}^i + (1 - 2\bar{\varepsilon}^i)\Phi(\lambda_i) \quad (38)$$

Here Φ is the cumulative distribution function of the standard normal distribution. According to formulas (27) and (28), we can obtain the following updating rules for a and b .

$$\Psi(a^{new}) - \Psi(a^{new} + b^{new}) \frac{\bar{\varepsilon}^i (1 - \Phi(\lambda_i))}{a^i [\bar{\varepsilon}^i + (1-2\bar{\varepsilon}^i)\Phi(\lambda_i)]} + \Psi(a^i) - \Psi(a^i + b^i + 1) \quad (39)$$

$$\Psi(b^{new}) - \Psi(a^{new} + b^{new}) \frac{\bar{\varepsilon}^i (1 - \Phi(\lambda_i))}{b^i [\bar{\varepsilon}^i + (1-2\bar{\varepsilon}^i)\Phi(\lambda_i)]} + \Psi(b^i) - \Psi(a^i + b^i + 1) \quad (40)$$

Where $\Psi(x) = d \log(\Gamma(x))$ and Γ is the gamma function. As for the fact that $\Psi(x)$ is a non-linear function, we can only use numerical solution to update a^{new} and b^{new} . In order to avoid the computational complexity, the expectation propagation of ε and ε^2 are used instead of the expectation propagation of $\log(\varepsilon)$ and $\log(1-\varepsilon)$. Although it is not guaranteed to minimize the KL divergence, the results are still accurate according to (Hernández-Lobato and Hernández-Lobato, 2008) and (Miguel Hernández-Lobato, et al., 2011). In other words, we can use (41) and (42) to update the value of a^{new} and b^{new}

$$\mathbb{E}_{\tilde{\varepsilon}^i}[\varepsilon] = \mathbb{E}_{t_i Q^i}[\varepsilon] \quad (41)$$

$$\mathbb{E}_{\tilde{\varepsilon}^i}[\varepsilon \circ \varepsilon] = \mathbb{E}_{t_i Q^i}[\varepsilon \circ \varepsilon] \quad (42)$$

After simplifying the equations we get

$$a^{new} = \frac{\mathbb{E}_{t_i Q^i}[\varepsilon] - \mathbb{E}_{t_i Q^i}[\varepsilon^2]}{\mathbb{E}_{t_i Q^i}[\varepsilon^2] - \mathbb{E}_{t_i Q^i}[\varepsilon]^2} \mathbb{E}_{t_i Q^i}[\varepsilon] \quad (43)$$

$$b^{new} = \frac{\mathbb{E}_{t_i Q^i}[\varepsilon] - \mathbb{E}_{t_i Q^i}[\varepsilon^2]}{\mathbb{E}_{t_i Q^i}[\varepsilon^2] - \mathbb{E}_{t_i Q^i}[\varepsilon]^2} (1 - \mathbb{E}_{t_i Q^i}[\varepsilon]) \quad (44)$$

In the above two equations, we have

$$\mathbb{E}_{t_i Q^i}[\varepsilon] = \frac{1}{Z_i(a^i + b^i + 1)} [\Phi(\lambda_i)(1 - \bar{\varepsilon}^i)a^i + (1 - \Phi(\lambda_i))\bar{\varepsilon}^i(a^i + 1)] \quad (45)$$

$$\mathbb{E}_{t_i Q^i}[\varepsilon^2] = \frac{a^i + 1}{Z_i(a^i + b^i + 1)(a^i + b^i + 2)} [\Phi(\lambda_i)(1 - \bar{\varepsilon}^i)a^i + (1 - \Phi(\lambda_i))\bar{\varepsilon}^i(a^i + 2)] \quad (46)$$

As for the approximation of t_{n+1} , or namely $p(\beta | z)$, we have (24), (25) and (26) here according to the infer from the minimum KL divergence between $t_i Q^i$ and $\tilde{t}_i Q^i$

$$\mathbb{E}_{\tilde{t}_i Q^i}[\beta] = \mathbb{E}_{t_i Q^i}[\beta] \quad (47)$$

$$\mathbb{E}_{\tilde{t}_i Q^i}[\beta \circ \beta] = \mathbb{E}_{t_i Q^i}[\beta \circ \beta] \quad (48)$$

$$\mathbb{E}_{\tilde{t}_i Q^i}[\mathbf{z}] = \mathbb{E}_{t_i Q^i}[\mathbf{z}] \quad (49)$$

Based on the above three equations, the rules for updating \mathbf{m} , \mathbf{v} and ρ can be derived as follows.

$$\mathbf{m}^{new} = \mathbf{m}^i + k' \circ \mathbf{v}^i \quad (50)$$

$$\mathbf{v}^{new} = \mathbf{v}^i - k''' \circ \mathbf{v}^i \circ \mathbf{v}^i \quad (51)$$

$$\rho^{new} = \rho^i + \rho^i(\rho^i)\nabla_{\rho} \log Z_i \quad (52)$$

$$\rho^{new} = \rho^i + \frac{(g'' - g''')\rho(1 - \rho^i)}{\rho^i g'' + (1 - \rho^i)g'''} \quad (53)$$

$$\rho^{new} = \rho^i \circ g'' \circ (\rho^i \circ g'' + (1 - \rho^i)g''') \quad (54)$$

k' , k''' , g'' and g''' in above equations can be given as follows

$$g'' = \mathcal{N}(0, \mathbf{m}^i, \mathbf{v}^i + \sigma^2) \quad (55)$$

$$g''' = \mathcal{N}(0, \mathbf{m}^i, \mathbf{v}^i) \quad (56)$$

$$g' = \rho^i \circ g'' + (1 - \rho^i) \circ g''' \quad (57)$$

$$k' = -\frac{\rho^i \circ g'' \circ \mathbf{m}^i}{g' \circ (\mathbf{v}^i + \sigma^2)} - \frac{(1 - \rho^i) \circ g''' \circ \mathbf{m}^i}{g' \circ \mathbf{v}^i} \quad (58)$$

$$k'' = \frac{\rho^i \circ g'' \circ \mathbf{m}^i \circ \mathbf{m}^i}{g' \circ (\mathbf{v}^i + \sigma^2) \circ (\mathbf{v}^i + \sigma^2)} - \frac{\rho^i \circ g''}{g' \circ (\mathbf{v}^i + \sigma^2)} + \frac{(1 - \rho^i) \circ g''' \circ \mathbf{m}^i \circ \mathbf{m}^i}{g' \circ \mathbf{v}^i \circ \mathbf{v}^i} - \frac{(1 - \rho^i) \circ g'''}{g' \circ \mathbf{v}^i} \quad (59)$$

$$k''' = k' \circ k' - k'' \quad (60)$$

Z_{n+1} can be given as follows while $Beta(\varepsilon|a^{\setminus i}, b^{\setminus i})$ does not contain β and Z_i

$$Z_i = \int (\mathcal{N}(\beta_i, 0, \sigma_i^2)^{z_i} \delta(\beta_i)^{(1-z_i)}) \prod_{j=1}^d \mathcal{N}(\beta_j | m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j | \rho_j^{\setminus i}) d\beta dz = \prod_{j=0}^d g_j' \quad (61)$$

As for the approximation of \tilde{t}_i for $t_i \in \mathcal{P}(\mathbf{z} | G, \lambda, \gamma)$ ($i = n+2 \dots, n+|E|+2$), we have formula (4) here

$$\begin{aligned} p(\mathbf{z} | G, \lambda, \gamma) &= \frac{1}{Z} \exp(cz_0) \exp(\lambda \sum_{i=1}^d z_i + \gamma \sum_{\{u,v\} \in E} \left(\frac{z_u}{\sqrt{d_u}} - \frac{z_v}{\sqrt{d_v}} \right)^2 w(u, v)) \\ &= \frac{1}{Z} \exp(cz_0 + \lambda \sum_{i=1}^d z_i) \exp(\gamma \sum_{\{u,v\} \in E} \left(\frac{z_u}{\sqrt{d_u}} - \frac{z_v}{\sqrt{d_v}} \right)^2 w(u, v)) \end{aligned} \quad (62)$$

Firstly, we need to approximate the priori sparse term $\exp(cz_0 + \lambda \sum_{i=1}^d z_i)$ and the following formula holds

$$\mathbb{E}_{\tilde{t}_i Q^{\setminus i}}[\mathbf{z}] = \mathbb{E}_{t_i Q^{\setminus i}}[\mathbf{z}] \quad (63)$$

And Z_i can be calculated by

$$Z_i = \int (\exp(h_i z_i) \text{Beta}(\varepsilon | a^{\setminus i}, b^{\setminus i})) \prod_{j=1}^d \mathcal{N}(\beta_j | m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j | \rho_j^{\setminus i}) d\beta dz \quad (64)$$

As in the above equation, $\mathbf{h} = (h_0, h_1, \dots, h_d)^T$ is a $d+1$ -dimension vector of which the first component is 0 while the others are λ , we can do the simplification as follows:

$$\begin{aligned} Z_i &= \exp(h_i) \rho^{\setminus i} \int \prod_{j=0}^d \mathcal{N}(\beta_j | m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j | \rho_j^{\setminus i}) d\beta + \exp(-h_i) (1 - \rho^{\setminus i}) \int \prod_{j=0}^d \mathcal{N}(\beta_j | m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j | \rho_j^{\setminus i}) d\beta \\ & \end{aligned} \quad (65)$$

$$Z_i = \prod_{j=0}^d [\rho_j^{\setminus i} \exp(h_i) + (1 - \rho_j^{\setminus i}) \exp(-h_i)] \quad (66)$$

According to (63), the updating rule for ρ is written as follows

$$\rho^{\text{new}} = \rho^{\setminus i} + \rho^{\setminus i} (1 - \rho^{\setminus i}) \nabla_{\rho} \log Z_i \quad (67)$$

Combined with (66), we have

$$\rho^{\text{new}} = \exp(\mathbf{h}) \circ \rho^{\setminus i} \circ (\exp(\mathbf{h}) \circ \rho^{\setminus i} + \mathbf{I}(1 - \rho^{\setminus i}))^{-1} \quad (68)$$

As for the approximation of \tilde{t}_i for i in range $(n+3, n+|E|+2)$

$$Z_i = \int \left(\exp \left(\gamma \left(\frac{z_u}{\sqrt{d_u}} - \frac{z_v}{\sqrt{d_v}} \right)^2 \right) \right) \text{Beta}(\varepsilon | a^{\setminus i}, b^{\setminus i}) \prod_{j=1}^d \mathcal{N}(\beta_j | m_j^{\setminus i}, v_j^{\setminus i}) \text{Bern}(z_j | \rho_j^{\setminus i}) d\beta d\varepsilon \quad (69)$$

Assume the A_i, B_i, C_i, D_i can be given as follows

$$A_i = \rho_u^{\setminus i} \rho_v^{\setminus i} \exp(\gamma (\frac{1}{\sqrt{d_u}} - \frac{1}{\sqrt{d_v}})^2) \quad (70)$$

$$B_i = \rho_u^{\setminus i} (1 - \rho_v^{\setminus i}) \exp(\frac{\gamma}{d_u}) \quad (71)$$

$$C_i = (1 - \rho_u^i) \rho_v^i \exp\left(\frac{Y}{d_v}\right) \quad (72)$$

$$D_i = (1 - \rho_u^i)(1 - \rho_v^i) \quad (73)$$

The updating rule of ρ can be obtained as follows

$$\rho_u^{new} = \frac{A_i + B_i}{A_i + B_i + C_i + D_i} \quad (74)$$

$$\rho_v^{new} = \frac{A_i + C_i}{A_i + B_i + C_i + D_i} \quad (75)$$

Lastly, as for the approximation of $t_{n+|E|+3}$

$$Z_i = B(a, b) B(a_0, b_0)^{-1} B(a^i, b^i)^{-1} \quad (76)$$

According to the rules of propagating the expectation of ε and ε^2 , we have.

$$a^{new} = a_0 + a^i - 1$$

$$b^{new} = b_0 + b^i - 1$$

(77)

We can get the following updating rules based on the above expectation propagation algorithm.

$$\begin{aligned} \tilde{v}_i^{new} &= (v^{-1} - (v^i)^{-1})^{-1} \\ \tilde{\mathbf{m}}_i^{new} &= \tilde{v}_i^{new} \circ v^{-1} \circ m - \tilde{v}_i^{new} \circ (v^i)^{-1} \circ m^i \\ \tilde{c}_i^{new} &= \rho \circ (\rho^i)^{-1} \\ \tilde{d}_i^{new} &= (1 - \rho) (1 - \rho^i)^{-1} \\ \tilde{a}_i^{new} &= a - a^i \\ \tilde{b}_i^{new} &= b - b^i \end{aligned} \quad (78)$$

\tilde{s}_i is a constant and its updating rule is derived from \tilde{t}_i . When $i = 1, 2, \dots, n$

$$\tilde{s}_i^{new} = Z_i \sqrt{\prod_{j=0}^d \frac{\tilde{v}_{ij}^{new} + v_j^i}{\tilde{v}_{ij}^{new}}} \exp\left(\frac{1}{2} \sum_{j=0}^d \frac{(\tilde{m}_{ij}^{new} - m_j^i)^2}{\tilde{v}_{ij}^{new} + v_j^i}\right) \frac{B(a^i, b^i)}{B(a, b)} \quad (79)$$

When $i = n + 2, \dots, n + |E| + 2$, the updating rule of \tilde{s}_i becomes

$$\tilde{s}_i^{new} = Z_i \prod_{j=0}^d \sqrt{\frac{\tilde{v}_{ij}^{new} + v_j^i}{\tilde{v}_{ij}^{new}}} \exp\left(\frac{1}{2} \frac{(k_j^i)^2}{k_j^{iii}}\right) \quad (80)$$

When $i = n + 2, \dots, n + |E| + 2$, the updating rule of \tilde{s}_i becomes

$$\tilde{s}_i^{new} = Z_i \quad (81)$$

When $i = n + |E| + 3$, the updating rule becomes

$$\tilde{s}_i^{new} = B(a_0, b_0)^{-1} \quad (82)$$

Once the expectation propagation algorithm converges, we can approximate it according to the following formula

$$p(\mathbf{y}|\mathbf{x}, G, \lambda, \gamma) \approx \int \sum_z \prod_{i=1}^{n+|E|+3} \tilde{t}_i(\beta, \varepsilon, z) d\beta d\varepsilon \approx \hat{Z}^{-1} C (2\Pi)^{\frac{d}{2}} \exp\left(\frac{D}{2}\right) B(A, B) [\prod_{i=1}^{n+|E|+3} \tilde{s}_i] [\prod_{j=0}^d \sqrt{v_j}] \quad (83)$$

Where as

$$A = \sum_{i=1}^{n+|E|+3} \tilde{a}_i + 1 \quad (84)$$

$$B = \sum_{i=1}^{n+|E|+3} \tilde{b}_i + 1 \quad (85)$$

$$C = \prod_{j=0}^d (\prod_{i=1}^{n+|E|+3} \tilde{c}_{ij} + \prod_{i=1}^{n+|E|+3} \tilde{d}_{ij}) \quad (86)$$

$$D = \mathbf{m}^T (\mathbf{v}^{-1} \circ \mathbf{m}) - \sum_{i=1}^{n+|E|+3} \tilde{\mathbf{m}}_i^T (\tilde{v}_i^{-1} \circ \tilde{\mathbf{m}}_i) \quad (87)$$

\hat{Z} is the approximation of Z in (3). Finally, we can predict the label of new samples according to the following formula

$$p(\mathbf{y}^{test} | \mathbf{x}^{test}, \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) \approx \int \int p(\mathbf{y}^{test} | \mathbf{x}^{test}, \beta, \varepsilon, G, \lambda, \gamma) \sum_z p(\beta, \varepsilon, z | \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) d\beta d\varepsilon = \int \int [\varepsilon + (1 - 2\varepsilon)\phi(\mathbf{y}^{test} \beta \mathbf{x}^{test})] \sum_z \mathcal{Q}(\beta, \varepsilon, z) d\beta d\varepsilon \quad (88)$$

According to (14)

$$p(\mathbf{y}^{test} | \mathbf{x}^{test}, \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) \approx \int \int [\varepsilon + (1 - 2\varepsilon)\phi(\mathbf{y}^{test} \beta \mathbf{x}^{test})] \sum_z \text{Beta}(\varepsilon | a, b) \prod_{j=0}^d \mathcal{N}(\beta_j | m_j, v_j) \text{Bern}(z_j | \rho_j) d\beta d\varepsilon \quad (89)$$

After simplification, we have

$$p(\mathbf{y}^{test} | \mathbf{x}^{test}, \mathbf{y}, \mathbf{X}, G, \lambda, \gamma) \approx \bar{\varepsilon} + (1 - 2\bar{\varepsilon}) \phi\left(\frac{\mathbf{y}^{test} \mathbf{m}^T \mathbf{x}^{test}}{\sqrt{(\mathbf{v} \circ \mathbf{x}^{test})^T \mathbf{x}^{test}}}\right) \quad (90)$$

Where as

$$\bar{\varepsilon} = \frac{a}{a+b} \quad (91)$$

2 Feature selection in NBSBM

A relevant score was defined by equation (10) to quantify the relevance of a feature to the classification results. We applied equation (10) on the first dataset to extract features that are most relevant to the prostate cancer cell responses to Dasatinib. Supplementary table 1 shows those top-25 relevant genes that ranked by the relevant score. Among the top-ranked genes, CTNNB1, FGFR4, GRK6 and PHB2 are oncogenes that have been reported to play important role in prostate cancer development and progression (FitzGerald, et al., 2009; Linch, et al., 2017; Nakai, et al., 2019; Yang, et al., 2018). Then we did canonical

pathway enrichment analysis, those significantly enriched pathways were listed out in Supplementary table 2. MHC class II antigen presentation, Integration of energy metabolism, MAPK family signaling cascades, RAF/MAP kinase cascade, FLT3 Signaling pathways are top-enriched signaling pathways that correlated with the prostate cancer cell responses to Dasatinib, which was also reported by the literature (da Silva, et al., 2013; Mukherjee, et al., 2011; Younger, et al., 2007).

Gene Entrez ID	Gene Symbol	Relevant Score
1499	CTNNB1	0.9999
51005	AMDHD2	0.9999
2264	FGFR4	0.9998
2870	GRK6	0.9913
11331	PHB2	0.9913
8504	PEX3	0.9913
8851	CDK5R1	0.9913
80700	UBXN6	0.9913
8078	USP5	0.9913
9409	PEX16	0.9913
22826	DNAJC8	0.9913
7317	UBA1	0.9913
55968	NSFL1C	0.9913
3053	SERPIND1	0.9913
57591	MRTFA	0.9913
10635	RAD51AP1	0.9913
8541	PPFIA3	0.9913
4601	MXI1	0.9913
55844	PPP2R2D	0.9913
5526	PPP2R5B	0.9913
51400	PPME1	0.9913
3009	H1-5	0.9913
9989	PPP4R1	0.9913
57718	PPP4R4	0.9913

Supplemental Table 1 Top-25 most predictive genes for classifying prostate cancer cell responses to Dasatinib. Oncogenes such as CTNNB1, FGFR4, GRK6 and PHB2 are top-ranked.

Enriched Pathways	p-value
MHC class II antigen presentation	2.74E-06
Integration of energy metabolism	0.002008
MAPK family signaling cascades	0.002191
RAF/MAP kinase cascade	0.006537
FLT3 Signaling	0.006955
MAPK1/MAPK3 signaling	0.009353
interleukin signaling	0.012823

Rho GTPase cycle	0.015637
Downstream TCR signaling	0.027074
Signaling by Receptor Tyrosine Kinases	0.028316
RHO GTPases Activate Formins	0.037878

Supplemental Table 2 The most enriched signaling pathways in those top-100 ranked genes that are most relevant to prostate cancer cell response to Dasatinib. P-value was estimated using the fisher's exact test.

References

- da Silva, H.B., *et al.* (2013) Dissecting Major Signaling Pathways throughout the Development of Prostate Cancer, *Prostate Cancer*, **2013**, 23.
- FitzGerald, L.M., *et al.* (2009) Association of FGFR4 genetic polymorphisms with prostate cancer risk and prognosis, *Prostate Cancer Prostatic Dis*, **12**, 192-197.
- Hernández-Lobato, D. and Hernández-Lobato, J.M. (2008) Bayes Machines for binary classification, *Pattern Recognition Letters*, **29**, 1466-1473.
- Linch, M., *et al.* (2017) Intratumoural evolutionary landscape of high-risk prostate cancer: the PROGENY study of genomic and immune parameters, *Annals of oncology : official journal of the European Society for Medical Oncology*, **28**, 2472-2480.
- Miguel Hernández-Lobato, J., Hernández-Lobato, D. and Suárez, A. (2011) Network-based sparse Bayesian classification, *Pattern Recognition*, **44**, 886-900.
- Mukherjee, R., *et al.* (2011) Upregulation of MAPK pathway is associated with survival in castrate-resistant prostate cancer, *British journal of cancer*, **104**, 1920-1928.
- Nakai, A., *et al.* (2019) The COMMD3/8 complex determines GRK6 specificity for chemoattractant receptors, *The Journal of Experimental Medicine*, **216**, 1630-1647.
- Yang, J., Li, B. and He, Q.-Y. (2018) Significance of prohibitin domain family in tumorigenesis and its implication in cancer diagnosis and treatment, *Cell Death Dis*, **9**, 580-580.
- Younger, A.R., *et al.* (2007) HLA class II antigen presentation by prostate cancer cells, *Prostate Cancer Prostatic Dis*, **11**, 334.