

Detecting DNA modifications from SMRT sequencing data by modeling sequence context dependence of polymerase kinetic

Zhixing Feng^{1,4}, Gang Fang², Jonas Korlach³, Tyson Clark³, Khai Luong³, Xuegong Zhang¹, Wing Wong^{4,*}, and Eric Schadt^{3,5,*}

¹Tsinghua National Laboratory for Information Science and Technology, and Department of Automation, Tsinghua University, Beijing 100084, China; ²Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN 55455; ³Pacific Biosciences, 1380 Willow Road, Menlo Park, CA 94025; ⁴Department of Statistics, Stanford University, Stanford, CA 94305; ⁵Department of Genetics and Genomics Sciences, Mount Sinai School of Medicine, New York, NY 10029; *Correspondence: whwong@stanford.edu, eric.schadt@mssm.edu

Text S1

Hyperparameters estimation for alternative model

In this section, we explain how to estimate hyperparameters by assuming control sample is available. As hierarchical model without control data is basically a special case of hierarchical model with control data, one can simply remove y_0 , μ_0 and σ_0^2 to get hyperparameter estimation when control sample is not available, and algorithm is unchanged. When alternative model is true (see section **Hierarchical model with control data** and **Hierarchical model without control data** in the main text), posterior distribution of μ_i and σ_i^2 , $i = 0, 1, \dots, m$, are[1]

$$p(\mu_i | \mathbf{y}_i, \sigma_i^2, \theta, \kappa) = N\left(\frac{\kappa}{\kappa + n_i}\theta + \frac{n_i}{\kappa + n_i}\bar{y}_i, \frac{\sigma_i^2}{\kappa + n_i}\right)$$
$$p(\sigma_i^2 | \mathbf{y}_i, v, \tau^2) = \text{scaled inverse} - \chi^2(v + n_i, \tilde{\sigma}_i^2)$$

where

$$\tilde{\sigma}_i^2 = \frac{1}{v + n_i} \left(v\tau^2 + (n_i - 1)s_i^2 + \frac{\kappa n_i}{\kappa + n_i} (\bar{y}_i - \theta)^2 \right) \quad (1)$$

$\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$ and $s_i = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$. Posterior distribution of σ_c^2 is[1]

$$p(\sigma_c^2 | \mathbf{y}_c, v, \tau^2) = \text{scaled inverse} - \chi^2(v + n_c, \tilde{\sigma}_c^2)$$

where

$$\tilde{\sigma}_c^2 = \frac{1}{v + n_c} \left(v\tau^2 + (n_i - 1)s_c^2 \right)$$

We used posterior expectations of μ_i and σ_i^2 as their estimations, which are

$$\hat{\mu}_i = E(\mu_i | \mathbf{y}_i, \theta, \kappa) = \frac{\kappa}{\kappa + n_i} \theta + \frac{n_i}{\kappa + n_i} \bar{y}_i \quad (2)$$

where $i = 0, 1, \dots, m$.

$$\hat{\sigma}_i^2 = E(\sigma_i^2 | \mathbf{y}_i, v, \tau^2) = \frac{n_i + v}{n_i + v - 2} \tilde{\sigma}_i^2$$

where $i = c, 0, 1, \dots, m$. We estimate hyperparameters $(\theta, \kappa, v, \tau^2, \mu_c)$ from the data by maximizing the marginal log-likelihood function, which is

$$\begin{aligned}
L(\mathbf{y}_c, \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_m; \theta, \kappa, v, \tau^2, \mu_c) &= \log(p(\mathbf{y}_c, \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_m | \theta, \kappa, v, \tau^2, \mu_c)) \\
&= \log(p(\mathbf{y}_c | \theta, \kappa, v, \tau^2, \mu_c)) \\
&\quad + \sum_{i=0}^m \log(p(\mathbf{y}_i | \theta, \kappa, v, \tau^2))
\end{aligned}$$

where

$$\begin{aligned}
\log(p(\mathbf{y}_c | \theta, \kappa, v, \tau^2, \mu_c)) &= \frac{v}{2} \log(v\tau^2) + \log(\Gamma(\frac{v+n_i}{2})) \\
&\quad - \frac{v+n_i}{2} \log(\sum_{j=1}^{n_c} (y_{cj} - \mu_c)^2 + v\tau^2) - \log(\Gamma(\frac{v}{2})) \\
&\quad - \frac{n_c}{2} \log(\pi) \\
\sum_{i=0}^m \log(p(\mathbf{y}_i | \theta, \kappa, v, \tau^2)) &= \sum_{i=0}^m [\frac{1}{2} \log(\kappa) + \log(\Gamma(\frac{v+n_i}{2})) + \frac{v}{2} \log(v\tau^2) \\
&\quad - \frac{1}{2} \log(\kappa + n_i) - \log(\Gamma(\frac{v}{2})) - \frac{v+n_i}{2} \log((v+n_i)\tilde{\sigma}_i^2) \\
&\quad - \frac{n_i}{2} \log(\pi)]
\end{aligned}$$

It is obvious that $L(\mathbf{y}_c, \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_m; \theta, \kappa, v, \tau^2, \mu_c)$ can be maximized by setting $\mu_c = \frac{1}{n_c} \sum_{j=1}^{n_c} y_{cj}$. However, it is difficult to get a close form of $(\theta, \kappa, v, \tau^2)$ and we therefore adopted a EM algorithm (Algorithm 1) to esti-

mate them numerically.

In the EM procedure, μ_i and σ_i^2 were regarded as missing data, and the log-likelihood function with the complete data was

$$\begin{aligned}
& l(\mathbf{y}, \mu, \sigma^2; \theta, \kappa, \tau^2, v) \\
&= \log(p(\mathbf{y}|\mu, \sigma^2)) + \log(p(\mu|\sigma^2, \theta, \kappa)) + \log(p(\sigma^2|v, \tau^2)) \\
&= \log(p(\mathbf{y}|\mu, \sigma^2)) + \sum_{i=0}^m \log(p(\mu_i|\sigma_i^2, \theta, \kappa)) + \sum_{i=c,0,1,\dots,m} \log(p(\sigma_i^2|\tau^2, v)) \quad (3)
\end{aligned}$$

where $\mathbf{y} = (\mathbf{y}_c, \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_m)$, $\mu = (\mu_0, \mu_1, \dots, \mu_m)$ and $\sigma^2 = (\sigma_c^2, \sigma_0^2, \sigma_1^2, \dots, \sigma_m^2)$.

Initial values $(\theta_0, \kappa_0, \tau_0^2, v_0)$ were assigned to $(\theta, \kappa, \tau^2, v)$, and in the t th step ($t \geq 1$) of EM algorithm, $(\theta, \kappa, \tau^2, v)$ were updated by the optimal value $(\theta_{opt}, \kappa_{opt}, \tau_{opt}^2, v_{opt})$ maximizing $E(l(\mathbf{y}, \mu, \sigma^2; \theta, \kappa, \tau^2, v) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$, where $E(\cdot)$ is expectation in term of (μ, σ^2) , i.e. $\theta_t = \theta_{opt}, \kappa_t = \kappa_{opt}, \tau_t = \tau_{opt}^2, v_t = v_{opt}$. This procedure was repeated until convergence.

In the t th step of the EM algorithm, $(\theta_{opt}, \kappa_{opt})$ and (τ_{opt}^2, v_{opt}) can be obtained by maximizing posterior expectation of the second term and the third term of equation (3), i.e. $\sum_{i=0}^m E(\log(p(\mu_i|\sigma_i^2, \theta, \kappa)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$ and $\sum_{i=c,0,1,\dots,m} E(\log(p(\sigma_i^2|\tau^2, v)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$, respectively.

$$\frac{{}^1\Delta l}{E(l(\mathbf{y}, \mu, \sigma^2; \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})} = \frac{E(l(\mathbf{y}, \mu, \sigma^2; \theta_{opt}, \kappa_{opt}, \tau_{opt}^2, v_{opt}) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})}{E(l(\mathbf{y}, \mu, \sigma^2; \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})} -$$

Algorithm 1 EM algorithm for hyperparameters estimation

Assign initial values to hyperparameters, $\theta = \theta_0$, $\kappa = \kappa_0$, $v = v_0$ and $\tau^2 = \tau_0^2$.

Set $t = 1$

while $\Delta l \leq 0.1$ ¹ **do**

E-step: Calculate conditional expectation of log likelihood function in terms of μ_i and σ_i^2 , i.e. $E(l(\mathbf{y}, \mu, \sigma^2; \theta, \kappa, \tau^2, v) \mid \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$.

M-step: Update $(\theta, \kappa, v, \tau^2)$ by setting $\theta_t = \theta_{opt}, \kappa_t = \kappa_{opt}, v_t = v_{opt}, \tau_t = \tau_{opt}^2$, where $(\theta_{opt}, \kappa_{opt}, v_{opt}, \tau_{opt}^2)$ are hyperparameters maximizing $E(l(\mathbf{y}, \mu, \sigma^2; \theta, \kappa, \tau^2, v) \mid \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$

Set $t = t + 1$

end while

Estimating θ and κ By taking posterior expectation of the second term of equation (3), we can get

$$\begin{aligned} & \sum_{i=0}^m E(\log(p(\mu_i | \sigma_i^2, \theta, \kappa)) \mid \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) \\ &= -\frac{\kappa}{2} \sum_{i=0}^m \left(\frac{1}{\kappa_{t-1} + n_i} + \frac{(\theta - \hat{\mu}_{i(t-1)})^2}{\tilde{\sigma}_{i(t-1)}^2} \right) - \frac{m+1}{2} \log(\kappa) + C \end{aligned}$$

where $\tilde{\sigma}_{i(t-1)}^2$ and $\hat{\mu}_{i(t-1)}$ are estimated $\tilde{\sigma}_i^2$ and μ_i given $(\theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$ (equation (1) and (2)), and C is a constant, which doesn't contain any hyperparameters. $\sum_{i=0}^m E(\log(p(\mu_i | \sigma_i^2, \theta, \kappa)) \mid \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$ can be

maximized by θ_{opt} and κ_{opt} , which are

$$\begin{aligned}\theta_{opt} &= \sum_{i=0}^m \frac{\hat{\mu}_{i(t-1)}}{\tilde{\sigma}_{i(t-1)}^2} / \sum_{i=0}^m \frac{1}{\tilde{\sigma}_{i(t-1)}^2} \\ \kappa_{opt} &= (m+1) / \sum_{i=0}^m \left(\frac{1}{\kappa_{t-1} + n_i} + \frac{(\theta_{opt} - \hat{\mu}_{i(t-1)})^2}{\tilde{\sigma}_{i(t-1)}^2} \right)\end{aligned}$$

Estimating τ^2 and v By taking posterior expectation of the third term of equation (3), we can get

$$\begin{aligned}& \sum_{i=c,0,1,\dots,m} E(\log(p(\sigma_i^2 | \tau^2, v)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) \\ &= \frac{(m+2)v}{2} \log\left(\frac{v\tau^2}{2}\right) - \left(\frac{v}{2} + 1\right) \sum_{i=c,0,1,\dots,m} E(\log(\sigma_i^2) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) - \\ & \quad \frac{\tau^2 v}{2} \sum_{i=c,0,1,\dots,m} E\left(\frac{1}{\sigma_i^2} | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}\right) - (m+2) \log\left(\Gamma\left(\frac{v}{2}\right)\right) \\ &= \frac{(m+2)v}{2} \log\left(\frac{v\tau^2}{2}\right) - \left(\frac{v}{2} + 1\right) \sum_{i=c,0,1,\dots,m} \left(\log\left(\frac{(v_{t-1} + n_i)\tilde{\sigma}_{i(t-1)}^2}{2}\right) - \psi\left(\frac{v_{t-1} + n_i}{2}\right)\right) - \\ & \quad \frac{\tau^2 v}{2} \sum_{i=c,0,1,\dots,m} \frac{1}{\tilde{\sigma}_{i(t-1)}^2} - (m+2) \log\left(\Gamma\left(\frac{v}{2}\right)\right)\end{aligned}$$

where $\Gamma(\cdot)$ is gamma function and $\psi(\cdot)$ is digamma function.

By setting

$$\begin{cases} \frac{\partial}{\partial \tau^2} \sum_{i=c,0,1,\dots,m} E(\log(p(\sigma_i^2|\tau^2, v)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) = 0 \\ \frac{\partial}{\partial v} \sum_{i=c,0,1,\dots,m} E(\log(p(\sigma_i^2|\tau^2, v)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) = 0 \end{cases}$$

we got

$$\begin{cases} \frac{(m+2)v}{2\tau^2} - \frac{v}{2} \sum_{i=c,0,1,\dots,m} \frac{1}{\tilde{\sigma}_{i(t-1)}^2} = 0 \\ (m+2) \left(\log\left(\frac{v}{2}\right) + \log(\tau^2) - \psi\left(\frac{v}{2}\right) \right) - \sum_{i=c,0,1,\dots,m} \left(\log\left(\frac{(v_{t-1}+n_i)\tilde{\sigma}_{i(t-1)}^2}{2}\right) - \psi\left(\frac{v_{t-1}+n_i}{2}\right) \right) = 0 \end{cases}$$

we got close form solution of the above equations by using approximation

of digamma function, which is $\psi\left(\frac{v}{2}\right) \approx \log\left(\frac{v}{2}\right) - \frac{1}{v} - \frac{1}{3v^2}$.

$$\begin{aligned} \tau_{opt}^2 &= \frac{m+2}{\sum_{i=c,0,1,\dots,m} \frac{1}{\tilde{\sigma}_{i(t-1)}^2}} \\ v_{opt} &= \frac{2}{3\left(\sqrt{1 + \frac{4}{3}T} - 1\right)} \end{aligned}$$

where $T = \frac{1}{m+2} \sum_{i=c,0,1,\dots,m} \left(\log\left(\frac{(v_{t-1}+n_i)\tilde{\sigma}_{i(t-1)}^2}{2}\right) - \psi\left(\frac{v_{t-1}+n_i}{2}\right) \right) - \log(\tau_{opt}^2)$.

Hyperparameters estimation for null model

For hierarchical model with control data, we denote pooled Box-Cox transformed IPD of native and control data as \mathbf{y}_p (i.e. $(y_{c1}, y_{c2}, \dots, y_{cn_c}, y_{01}, y_{02}, \dots, y_{0n_0})$).

For hierarchical model without control data, we simply let \mathbf{y}_c , Box-Cox transformed IPD of native sample, equal to \mathbf{y}_p , because it is a special case of

hierarchical model with control data, in which y_0 is empty. We assume \mathbf{y}_p follows a normal distribution

$$\mathbf{y}_p \sim N(\mu_p, \sigma_p^2)$$

and $(\mu_p, \sigma_p^2), (\mu_1, \sigma_1^2), (\mu_2, \sigma_2^2), \dots, (\mu_m, \sigma_m^2)$ have the same prior distribution, which is

$$p(\sigma_i^2 | v, \tau^2) = \text{scaled inverse} - \chi^2(v, \tau^2)$$

$$p(\mu_i | \sigma_i^2, \theta, \kappa) = N\left(\theta, \frac{\sigma_i^2}{\kappa}\right)$$

where, $i = p, 1, \dots, m$. posterior distribution of μ_i and σ_i^2 , $i = p, 1, \dots, m$, are

$$p(\mu_i | \mathbf{y}_i, \sigma_i^2, \theta, \kappa) = N\left(\frac{\kappa}{\kappa + n_i}\theta + \frac{n_i}{\kappa + n_i}\bar{y}_i, \frac{\sigma_i^2}{\kappa + n_i}\right)$$

$$p(\sigma_i^2 | \mathbf{y}_i, v, \tau^2) = \text{scaled inverse} - \chi^2(v + n_i, \tilde{\sigma}_i^2)$$

where

$$\tilde{\sigma}_i^2 = \frac{1}{v + n_i} \left(v\tau^2 + (n_i - 1)s_i^2 + \frac{\kappa n_i}{\kappa + n_i} (\bar{y}_i - \theta)^2 \right)$$

[1]. We used posterior expectations of μ_i and σ_i^2 as their estimations, which are

$$\begin{aligned}\hat{\mu}_i &= E(\mu_i | \mathbf{y}_i, \theta, \kappa) = \frac{\kappa}{\kappa + n_i} \theta + \frac{n_i}{\kappa + n_i} \bar{y}_i \\ \hat{\sigma}_i^2 &= E(\sigma_i^2 | \mathbf{y}_i, v, \tau^2) = \frac{n_i + v}{n_i + v - 2} \tilde{\sigma}_i^2\end{aligned}$$

where $i = p, 1, \dots, m$. Like the previous section, we adopt a EM algorithm (Algorithm 1) to maximize marginal log-likelihood function $L(\mathbf{y}_p, \mathbf{y}_1, \dots, \mathbf{y}_m; \theta, \kappa, v, \tau^2)$.

In the EM procedure, μ_i and σ_i^2 were regarded as missing data, and the log-likelihood function with the complete data was

$$\begin{aligned}l(\mathbf{y}, \mu, \sigma^2; \theta, \kappa, \tau^2, v) \\ &= \log(p(\mathbf{y} | \mu, \sigma^2)) + \sum_{i=p,1,\dots,m} \log(p(\mu_i | \sigma_i^2, \theta, \kappa)) \\ &+ \sum_{i=p,1,\dots,m} \log(p(\sigma_i^2 | \tau^2, v))\end{aligned}\tag{4}$$

where $\mathbf{y} = (\mathbf{y}_p, \mathbf{y}_1, \dots, \mathbf{y}_m)$, $\mu = (\mu_p, \mu_1, \dots, \mu_m)$ and $\sigma^2 = (\sigma_p^2, \sigma_1^2, \dots, \sigma_m^2)$.

Estimating θ and κ By taking posterior expectation of the second term of equation (4), we can get

$$\begin{aligned} & \sum_{i=p,1,\dots,m} E(\log(p(\mu_i|\sigma_i^2, \theta, \kappa)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) \\ = & -\frac{\kappa}{2} \sum_{i=p,1,\dots,m} \left(\frac{1}{\kappa_{t-1} + n_i} + \frac{(\theta - \hat{\mu}_{i(t-1)})^2}{\tilde{\sigma}_{i(t-1)}^2} \right) - \frac{m+1}{2} \log(\kappa) + C \end{aligned}$$

Thus, $\sum_{i=p,1,\dots,m} E(\log(p(\mu_i|\sigma_i^2, \theta, \kappa)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1})$ can be maximized by

$$\begin{aligned} \theta_{opt} &= \sum_{i=p,1,\dots,m} \frac{\hat{\mu}_{i(t-1)}}{\tilde{\sigma}_{i(t-1)}^2} / \sum_{i=0}^m \frac{1}{\tilde{\sigma}_{i(t-1)}^2} \\ \kappa_{opt} &= (m+1) / \sum_{i=p,1,\dots,m} \left(\frac{1}{\kappa_{t-1} + n_i} + \frac{(\theta_{opt} - \hat{\mu}_{i(t-1)})^2}{\tilde{\sigma}_{i(t-1)}^2} \right) \end{aligned}$$

Estimating τ^2 and v By taking posterior expectation of the third term of equation (4), we can get

$$\begin{aligned} & \sum_{i=p,1,\dots,m} E(\log(p(\sigma_i^2|\tau^2, v)) | \mathbf{y}, \theta_{t-1}, \kappa_{t-1}, \tau_{t-1}^2, v_{t-1}) \\ = & \frac{(m+1)v}{2} \log\left(\frac{v\tau^2}{2}\right) - \left(\frac{v}{2} + 1\right) \sum_{i=p,1,\dots,m} \left(\log\left(\frac{(v_{t-1} + n_i)\tilde{\sigma}_{i(t-1)}^2}{2}\right) - \psi\left(\frac{v_{t-1} + n_i}{2}\right) \right) - \\ & \frac{\tau^2 v}{2} \sum_{i=p,1,\dots,m} \frac{1}{\tilde{\sigma}_{i(t-1)}^2} - (m+1) \log\left(\Gamma\left(\frac{v}{2}\right)\right) \end{aligned} \quad (5)$$

By using approximation of digamma function, which is $\psi(\frac{v}{2}) \approx \log(\frac{v}{2}) - \frac{1}{v} - \frac{1}{3v^2}$, (5) can be maximized by

$$\begin{aligned}\tau_{opt}^2 &= \frac{m+1}{\sum_{i=p,1,\dots,m} \frac{1}{\tilde{\sigma}_{i(t-1)}^2}} \\ v_{opt} &= \frac{2}{3(\sqrt{1 + \frac{4}{3}T} - 1)}\end{aligned}$$

where $T = \frac{1}{m+1} \sum_{i=p,1,\dots,m} \left(\log\left(\frac{(v_{t-1}+n_i)\tilde{\sigma}_{i(t-1)}^2}{2}\right) - \psi\left(\frac{v_{t-1}+n_i}{2}\right) \right) - \log(\tau_{opt}^2)$.

References

- [1] Gelman A, Carlin JB, Stern HS, Rubin DB (2003) Bayesian Data Analysis. Chapman and Hall/CRC, 2nd edition.