

1. ISLET model

(1) notation:

gene G
 cell type K
 subject J (NU in code) $Q = JK$
 length $N = \sum_{j=1}^J T_j$ (NS in code, # of samples)

★ If we consider age effect, $X \in [0, 1]$ $N \times 4K$

$N \times G$ $N \times 2K$ (df for slope) $N \times Q$ ($Q = JK$)

$$Y = X\beta + Au + \varepsilon$$

$u =$ cell type 1 for subject 1 ~ 50

gene 1 ... gene G

u_{11}
 u_{21}
 \vdots
 u_{J1}
 u_{12}
 \vdots
 u_{JK}

$$X = \begin{pmatrix} \theta_{111} & \theta_{112} & \dots & \theta_{11K} & z_1\theta_{111} & z_1\theta_{112} & \dots & z_1\theta_{11K} \\ \theta_{121} & \theta_{122} & \dots & \theta_{12K} & z_1\theta_{121} & z_1\theta_{122} & \dots & z_1\theta_{12K} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \theta_{1T_11} & \theta_{1T_12} & \dots & \theta_{1T_1K} & z_1\theta_{1T_11} & z_1\theta_{1T_12} & \dots & z_1\theta_{1T_1K} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \theta_{J11} & \theta_{J12} & \dots & \theta_{J1K} & z_J\theta_{J11} & z_J\theta_{J12} & \dots & z_J\theta_{J1K} \\ \theta_{J21} & \theta_{J22} & \dots & \theta_{J2K} & z_J\theta_{J21} & z_J\theta_{J22} & \dots & z_J\theta_{J2K} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \theta_{JT_11} & \theta_{JT_12} & \dots & \theta_{JT_1K} & z_J\theta_{JT_11} & z_J\theta_{JT_12} & \dots & z_J\theta_{JT_1K} \end{pmatrix}_{N \times 2K}$$

$$A_{11} = \begin{bmatrix} \theta_{111} \\ \theta_{121} \\ \vdots \\ \theta_{1T_11} \end{bmatrix}$$

$N = NS = NU \times T$ (5) 第1个 subject 在第1个 cell type
 different proportion across time

$$A = \begin{pmatrix} a_{11} & 0 & 0 & 0 & \dots & a_{1K} & 0 & 0 & 0 \\ 0 & a_{21} & 0 & 0 & \dots & 0 & a_{2K} & 0 & 0 \\ 0 & 0 & \ddots & 0 & \dots & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & a_{J1} & \dots & 0 & 0 & 0 & a_{JK} \end{pmatrix}_{N \times Q}$$

where $a_{jk} := (\theta_{j1k}, \theta_{j2k}, \dots, \theta_{jT_kk})'$ is simply a reorganized vector of cell type proportions, $Q = NU \times K$ to align with random effect u .

(2) assumption:

① $\varepsilon \sim \text{Normal}(0, \sigma_0^2)$

for every gene g and every subject j

② $u \sim \text{MVN}(\vec{0}_{Q \times 1}, \Sigma_u)$

for every gene g .

with $\Sigma_u = \text{Cov}(u) = \begin{pmatrix} \sigma_1^2 I_J & & 0 \\ & \sigma_1^2 I_J & \dots \\ 0 & \dots & \sigma_K^2 I_J \end{pmatrix}$

③ $y|u \sim N(X\beta + Au, \sigma_0^2 I)$ - $W := (y, u) = (W_{\text{observed}}, W_{\text{missed}})$

$\Rightarrow \begin{pmatrix} y \\ u \end{pmatrix} \sim N\left(\begin{bmatrix} X\beta \\ 0 \end{bmatrix}, \begin{bmatrix} A\Sigma_u A^T & A\Sigma_u \\ \Sigma_u^T A^T & \Sigma_u \end{bmatrix}\right)$ let $\eta = \frac{\text{baseline of ctrl group}}{\sigma_1^2, \dots, \sigma_K^2, \sigma_0^2} \beta_1, \dots, \beta_K$ group effect

If we want to test age,

$\eta = (m_1, \dots, m_K, \beta_1, \dots, \beta_K, \alpha_1, \dots, \alpha_K, \delta_1, \dots, \delta_K, \sigma_1^2, \dots, \sigma_K^2, \sigma_0^2)$
 age effect (across all sample) age + group effect

a_{11} $k=1$ (第1个 cell type, 5个时间点) a_{22} a_{12}

	V1	V2	V3	V4	...	V50
1	0.02230079	0.00000000	0.00000000	0.00000000		
2	0.01501410	0.00000000	0.00000000	0.00000000		
3	0.01314183	0.00000000	0.00000000	0.00000000		
4	0.04691058	0.00000000	0.00000000	0.00000000		
5	0.04692466	0.00000000	0.00000000	0.00000000		
6	0.00000000	0.07183477	0.00000000	0.00000000		
7	0.00000000	0.01136901	0.00000000	0.00000000		
8	0.00000000	0.06145904	0.00000000	0.00000000		
9	0.00000000	0.02264138	0.00000000	0.00000000		
10	0.00000000	0.04052955	0.00000000	0.00000000		
11	0.00000000	0.00000000	0.03578315	0.00000000		
12	0.00000000	0.00000000	0.08965811	0.00000000		
13	0.00000000	0.00000000	0.05633649	0.00000000		
14	0.00000000	0.00000000	0.08530417	0.00000000		
15	0.00000000	0.00000000	0.01453683	0.00000000		
16	0.00000000	0.00000000	0.00000000	0.05460231		

$u \in \mathbb{R}^{Q \times 1}$
 $Q = 50 \times 6 = 300$ in our case
 $NS = 250$
 $NU = 50$
 $T = 5$

Note: In islet code, the output will be a $G \times K$ matrix for each subject j .
(since slope are current unavailable), with each element m_{gk} for ctrl
group. $\beta_{gk} + m_{gk}$ for case group for gene g , cell type k .

joint (log-likelihood)

$$l(\eta; y, u) = \log[f(y|u; X, A, \beta, \sigma) \cdot f(u; X, A, \beta, \sigma)]$$

$$\propto \frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - X\beta - Au)^T (y - X\beta - Au) - \frac{J}{2} \sum_{k=1}^K \log(\sigma_k^2) - \frac{1}{2} \sum_{k=1}^K \frac{1}{\sigma_k^2} u_k^T u_k$$

Details: check ISLET Additional file 5. Appendix

parameters of interest: $\sigma^2, \sigma_1^2, \sigma_2^2, \dots, \sigma_K^2, \beta, u$ random effect
solved by EM algorithm fixed effect

The EM algorithm calculation will then follow naturally.

E-step:

$$E[u|w_{obs} = y] = \Sigma_u A' V^{-1} (y - X\beta) \quad V = A \Sigma_u A' + \sigma_0^2 I_N$$

$$E[s's|w_{obs} = y] = tr(A \Sigma_p A') + (A \mu_p + X\beta - y)' (A \mu_p + X\beta - y)$$

$$E[u'_k u_k | w_{obs} = y] = tr(\Sigma_{p_k}) + \mu'_{p_k} \mu_{p_k}$$

Here, $s = Au + X\beta - y$, $V := A \Sigma_u A' + \sigma_0^2 I$, Σ_{p_k} is the k th diagonal block of matrix Σ_p , and μ_{p_k} is the k th sub-vector in μ_p .

M-step:

For the $(t+1)^{th}$ iteration given the t^{th} iteration:

$$\hat{\beta}^{(t+1)} = (X'X)^{-1} X' (y - A E_{\eta^{(t)}}(u^{(t)}))$$

$$\hat{\sigma}_0^{2(t+1)} = \frac{E_{\eta^{(t)}}[s's|w_{obs} = y]}{N}$$

$$\hat{\sigma}_k^{2(t+1)} = \frac{E_{\eta^{(t)}}[u'_k u_k | w_{obs} = y]}{J}$$

Lemma 2.1. If $X = (X_1, X_2)$, and $X \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}\right)$, then $[X_1|X_2] \sim N(\mu_{1|2}, \Sigma_{1|2})$, where $\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(X_2 - \mu_2)$ and $\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$.

Then we have

$$[w_{mis}|w_{obs}] = [u|y] \sim N_Q(\mu_p, \Sigma_p)$$

and

$$\mu_p = \Sigma_u A' V^{-1} (y - X\beta) \quad \Sigma_p = \Sigma_u - \Sigma_u A' V^{-1} A \Sigma_u$$

Now we let $s = Au + X\beta - y$ and $V := A \Sigma_u A' + \sigma_0^2 I_N$. To apply Expectation-Maximization (EM) algorithm, we need to evaluate $E_{u|y, \eta}[l(\eta; y, u)]$, and consequently $E[s's|y, \eta]$, $E[u'_k u_k | Y, \eta]$.

模型现存的问题: (1) 不收敛. source code 只设置了14次迭代.

(2) $u=0$ 的情况解不出来 (但这种情况 real life 少见)

simulation 中, $u=0$ 时. 变为简单的解矩阵问题: $y = X\beta + \varepsilon$
with known X, y , want to solve β .

考虑 update:

(1) 迭代. 代码中列出了2个可能作为 stopping criteria 的量

B_change_val 和 $diff2$.

$$B_change_val := \text{sum}(|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}|) / \text{length}(\hat{\beta}^{(t)})$$

$$diff2 := \text{sum}(|E u^{(t+1)} - E u^{(t)}|) / (\text{length}(E u^{(t)}) \cdot \bar{y}^2)$$

(2) 代码效率: 尤其涉及 likelihood ratio test. 代码慢.

(3) 解出 time-wise 的 individual reference panel (4) 其他算法?

$$\bar{y} = \text{mean}(\text{colMeans}(y))$$

```
B_sum_val= 2046.932
B_change_val= 6.555662
B_change_prop= 0.3202677 %
Random effect diff2= 2.938226e-06
B_sum_val= 2045.422
B_change_val= 6.16733
B_change_prop= 0.3015187 %
Random effect diff2= 2.80272e-06
B_sum_val= 2044.006
B_change_val= 5.819443
B_change_prop= 0.2847077 %
Random effect diff2= 2.679876e-06
```

2. dispersion - estimation

① mean, variance, dispersion 三者的关系?

real life 中 variance 和 dispersion 的取值?

DESeq2:

dispersion max: 24

mean (raw count): 80

< mean (log scale): 3

variance (raw count): $= 6 \times 10^8$

(with mean = 3000)

对于 negative ~ binomial GLM:

$$y \sim \text{NB}(\mu, \phi)$$

$$Y_i | \lambda_i \sim \text{Poisson}(\lambda_i)$$

$$\lambda_i \sim \text{Gamma}(\mu_i, \phi)$$

It is straightforward to show, under the hierarchical model, that

$$E[Y_i] = \mu_i, \quad \text{var}[Y_i] = \mu_i + \phi \mu_i^2$$

where the variance contains an overdispersion term $\phi \mu_i^2$. The larger ϕ , the greater the overdispersion. The PMF of this NB is:

$$P(y_i; \mu_i, \alpha) = \frac{\Gamma(y_i + \alpha)}{\Gamma(y_i + 1)\Gamma(\alpha)} \left(\frac{\mu_i}{\mu_i + \alpha} \right)^{y_i} \left(1 - \frac{\mu_i}{\mu_i + \alpha} \right)^\alpha.$$

where $\alpha = 1/\phi$ and $\Gamma()$ is the gamma function, so that $\text{var}[Y_i] = \mu_i + \mu_i^2/\alpha$.

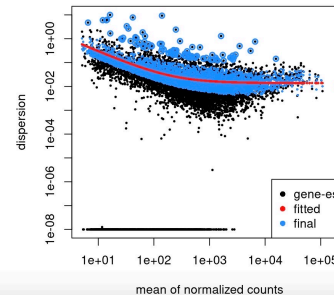
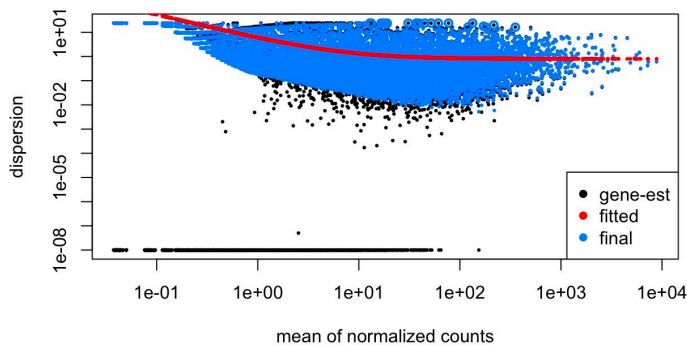
Negative binomial GLMs give larger standard errors than the corresponding Poisson GLMs, depending on the size of $k = 1/\psi$. On the other hand, the coefficient estimates $\hat{\beta}_j$ from a negative binomial GLM may be similar to those produced from the corresponding Poisson GLM. The negative binomial GLM gives less weight to observations with large μ_i than does the Poisson GLM, and relatively more weight to observations with small μ_i , so the coefficients

② dispersion shrinkage?

Dispersion plot and fitting alternatives

Plotting the dispersion estimates is a useful diagnostic. The dispersion plot below is typical, with the final estimates shrunk from the gene-wise estimates towards the fitted estimates. Some gene-wise estimates are flagged as outliers and not shrunk towards the fitted value, (this outlier detection is described in the manual page for `estimateDispersionMAF`). The amount of shrinkage can be more or less than seen here, depending on the sample size, the number of coefficients, the row mean and the variability of the gene-wise estimates.

```
plotDispEsts(dds)
```



写 writing sample. 11月 URA 读期 + 提交申请 (帮我改)

后读吧?