

varianceExplainTheory

December 28, 2017

A linear mixed model in model organism association mapping is typically expressed as $y = x + Zu + e$ where y is an $n \times 1$ vector of observed phenotypes, and x is an $n \times q$ matrix of fixed effects including mean, SNPs, and other confounding variables. u is a $q \times 1$ vector representing coefficients of the fixed effects. Z is an $n \times t$ incidence matrix mapping each observed phenotype to one of t inbred strains. And in general Z is None, and will be ignored at some positions in the following document. u is the random effect of the mixed model with $Var(u) = \frac{1}{g}K$, where K is the $t \times t$ kinship matrix, and e is an $n \times n$ matrix of residual effect such that $Var(e) = \frac{1}{e}$. The overall phenotypic variance-covariance matrix can be represented as $V = \frac{1}{g}K + \frac{1}{e}I$. This part, in italic font, explains likelihood distribution and the symbols have no relationship with context. For linear model $y = x$. The distribution of $y \sim N(x, \frac{1}{g}I)$. So the likelihood $L(y; \frac{1}{g}) = \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{(y-x)(y-x)}{2\frac{1}{g}}}$

For the linear mixed model: $y \sim N(x, \frac{1}{g}K + \frac{1}{e}I) = N(x, \frac{1}{g}H)$ $H = \frac{1}{g}V = K + I$ (there is a typo in the emma paper) $\frac{1}{g} likelihood(y; \frac{1}{g}) = \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{(y-x)H^{-1}(y-x)}{2\frac{1}{g}}}$ ML (maximum likelihood): $l_F(y; \frac{1}{g}) = \frac{1}{2}(-n\log(2\frac{1}{g}) - \log|H| - \frac{1}{2}(y-x)H^{-1}(y-x))$ REML: $l_R(y; \frac{1}{g}) = l_F(y; \frac{1}{g}) + \frac{1}{2}(q\log(2\frac{1}{g}) + \log|xx| - \log|xH^{-1}x|)$ Gradient of the LMM log likelihood w.r.t. $l_R(y; \frac{1}{g}) = \frac{d - \frac{1}{2g}(y-x)^T(K+I)^{-1}(y-x)}{d} = \frac{1}{2g}(-x^T(k+I)^{-1}y + x^T(k+I)^{-1}x)$ set gradient to zero: $x^TH^{-1}x = x^TH^{-1}y$ could be estimated as $\hat{x} = (xH^{-1}x)^{-1}xH^{-1}y$

Note that this solution is analogous to the ML solution of the linear regression $(x^Tx)^{-1}x^Ty$

For ML: $\frac{1}{g} = \frac{R}{n}$ Then: $l_F(y; \hat{g}, \hat{x}) = \frac{1}{2}(-n\log\frac{2R}{n} - \log|H| - n)$

For REML: $\frac{1}{g} = \frac{R}{n-q}$ $R = (y-x)H^{-1}(y-x)$ $H = K + I = U_F(i + \dots, n +)U_F$

$K = U_F U_F$ with eigen decomposition so $\log|H| = \sum_{i=1}^n \log(i +)$ And $R = (y-x)H^{-1}(y-x) = y(I - x(xH^{-1}x)^{-1}xH^{-1})H^{-1}(I - x(xH^{-1}x)^{-1}xH^{-1})y = yPH^{-1}Py$ P is defined as $P = I - x(xH^{-1}x)^{-1}xH^{-1}$ $S = I - X(XX^T)^{-1}X^T$ $SHS = S(K+I)S$ do eigen decomposition And $(SHS)(PH^{-1}P)(SHS) = SHS (PH^{-1}P)(SHS)(PH^{-1}P) = (PH^{-1}P)PS = P$ (there is a typo in the emma paper) and $SP = S$ $SHS = [U_R, W_R]diag(1 + \dots, n-q + , 0, \dots, 0)[U_R, W_R]$
 $= U_R diag(1 + \dots, n-q, 0, \dots, 0)U_R$

U_R is an $n \times (n - q)$ eigenvector matrix corresponding to the nonzero eigenvalues. W_R is an $n \times q$ eigenvector matrix corresponding to zero eigenvalues. Here the eigen decomposition of H and SHS do not dependent on any unknown parameter and could be done directly

so $PH^{-1}P = (SHS)^+ = U_R diag((s +)^{-1})U_R$ here $(.)^+$ denotes the pseudo-inverse of a matrix

Let $U_R y = [1, 2, \dots, n_q]$ and $R = yPH^{-1}Py = (U_R y)diag((s +)^{-1})(U_R y) = \sum_{s=1}^{n-q} \frac{2}{s + }$

Then for ML: $l_F(y; \hat{g}, \hat{x}) = \frac{1}{2}n\log\frac{n}{2} - n - n\log\sum_{s=1}^{n-q} \frac{2}{s + } - \sum_{i=1}^n \log(i +)$

$(SHS)(SHS)^+ = (SHS)(PH^{-1}P) = SHPH^{-1}P = SP = S$ On the other hand $(SHS)(SHS)^+ = (U_R \text{diag}(s +)U_R)(U_R \text{diag}((s +)^{-1})U_R) = U_R U_R$ so $U_R U_R = S = I$ taking account $\frac{2}{s} = \frac{R}{n-q}$
 Then for REML: $l_R(y; \hat{g}, \hat{\epsilon}) = \frac{1}{2}(n - q) \log \frac{n-q}{2} - (n - q) - (n - q) \log \sum_{s=1}^{n-q} \frac{2}{s} - \sum_{s=1}^{n-q} \log(s +)$ The
 derivatives of these functions: ML: $f_F = \frac{n}{2} \frac{\sum_s \frac{2}{(s+)^2}}{\frac{2}{s+}} - \frac{1}{2} \sum_i \frac{1}{i+}$ REML: $f_F = \frac{n-q}{2} \frac{\sum_s \frac{2}{(s+)^2}}{\frac{2}{s+}} - \frac{1}{2} \sum_i \frac{1}{i+}$ if
 we could find B such that $BB = H = \frac{V}{s} = K + I$ we can substitute $y^* = B^{-1}y$, $x^* = B^{-1}x$ and
 $* = B^{-1}(Zu +)$ (now * includes both random effects and errors) to get $y^* = x^* + *$ $Var(*) =$
 $Var(B^{-1}(Zu +)) = B^{-1}V(B^{-1}) = \frac{2}{s} B^{-1}H(B^{-1}) = \frac{2}{s} B^{-1}BB(B^{-1}) = \frac{2}{s} I$ The value of the residual
 sum of squares (RSS) from solving the transformed equation $y^* = X^* \beta + \epsilon^*$ is the Mahalanobis
 RSS for the original equation $y = x\beta + Zu + \epsilon$.

Taking advantage of the eigen decomposition of H performed in the EMMA algorithm, the
 computation of a valid B^{-1} can be simplified to $B^{-1} = \text{diag}(1/\sqrt{\xi_1 + \delta}, \dots, 1/\sqrt{\xi_n + \delta}) U_F' B^{-1}$ is
 H_sqrt_inv in the code

could be estimate by solving $y^* = x^* + *$

The F test here is performed on y^* and x^* In the code h0_rss is the Residual sum of squares
 without the effecte of fixed variable under testing (H0) mahalanobis_rss is the Residual sum of
 squares with the effecte of fixed variable under testing (H1)

If here x is the intercept and all the significant genotypic variants, after being calculated Then
 residuals = y - x Variance explained by significant genotypic variants is (var(y)-residuals.T * resid-
 uals) / var(y)