

利用回交或 $F_{2:3}$ 家系世代鉴定数量性状 主基因-多基因混合遗传模型*

盖钧镒** 王建康***

(南京农业大学大豆研究所, 江苏南京, 210095)

提 要 主-多基因性状在分离世代中表现出多峰或偏态现象, 具有混合分布的特征。本文利用混合分布理论提出从回交或 $F_{2:3}$ 家系个别分离世代鉴定主基因的存在并估计其遗传效应的统计分析方法。利用所提出的方法可以鉴别主基因的存在并估计其遗传效应和方差等遗传参数。

关键词 数量性状; 主基因和多基因混合遗传; 混合分布; EM 算法

大量试验数据和 QTL (quantitative trait loci) 作图结果表明: 控制数量性状的基因体系中既有遗传效应较大的主基因, 又有遗传效应较小的微效多基因, 称之为主基因-多基因混合遗传^[3, 5, 6, 11, 13, 16]。Elkind 和 Cahaner^[8, 9]提出单基因-多基因 (single gene-polygene, SG-PG) 遗传模型, 其单基因是可鉴定的。莫惠栋^[3]分析了一对主基因存在时, 主基因-多基因混合遗传性状在各个世代的遗传组成以及遗传参数的估计问题, 并把这类性状称为质量-数量性状, 由于 F_2 代的分组趋势不明显, 作者建议采用后裔测验的方法, 然后通过聚类分析确定 F_2 个体的主基因基因型。Loisel^[13]研究了在 F_2 世代中探测主基因存在的似然比统计量的渐近性质, 由此可对主基因是否存在、主基因是否表现为加性或完全显性进行检验。姜长鉴^[1]将 Loisel^[13]的结果用于大麦矮秆突变体与正常秆品系间杂交产生的 F_2 代株高性状的遗传分析。王建康和盖钧镒^[5, 6]提出利用混合分布理论在 F_2 群体中鉴定数量性状主基因的存在和估计遗传参数的方法。这些结果虽然可以检验主基因的加显性并估计主基因的加显性效应和主基因遗传方差, 但 F_2 世代受环境影响较大, 检测主基因的功效不高, 对多基因的存在无法作进一步鉴定。本文提出利用混合分布理论在回交和 $F_{2:3}$ 家系世代中鉴定主基因的存在并估计其遗传效应的统计分析方法。

1 主基因-多基因混合遗传模型

设来自杂种分离群体个体的表现型可分解为:

$$x = m + g + c + e$$

其中 m 为群体平均值, g 为主基因效应值, c 为多基因效应值, e 为环境效应值。此处假定主

* 国家 863 项目部分内容和河南省科委攻关项目

致谢: 南京农业大学农学系朱立宏教授惠允使用水稻南京 6 号与广丛回交世代株高分布的资料, 谨致谢忱。

** 联系作者, E-mail: sri@njau.edu.cn

*** 现在河南省农业科学院科学实验中心工作, 河南郑州, 450002

收稿日期: 1997-03-12, 收到修改稿日期: 1997-10-29

基因与多基因之间不存在互作, 基因型效应与环境效应是相互独立的。主基因效应 g 为固定值, 主基因方差用 σ_{mg}^2 表示; 多基因效应 $c \sim N(0, \sigma_{pc}^2)$ 是随机变量, σ_{pc}^2 为多基因效应值的方差; 环境效应 $e \sim N(0, \sigma_e^2)$ 是随机变量, σ_e^2 为环境方差, σ_p^2 表示总表型方差, 这样遗传率可分解为主基因遗传率 $h_{mg}^2 = \sigma_{mg}^2 / \sigma_p^2$ 及多基因遗传率 $h_{pc}^2 = \sigma_{pc}^2 / \sigma_p^2$ 。如果主基因有 k 个不同的基因型值, 那么分离群体表现为 k 个正态分布的混合, 通过对混合群体的分解可将主基因的效应分解出来。除主基因之外的变异是多基因变异和环境变异的混合, 在只有个别分离群体的情况下, 不能将二者分解开来, 记 $\sigma^2 = \sigma_{pg}^2 + \sigma_e^2$ 。

2 混合分布分析方法

混合分布理论在实际中有广泛的应用^[4, 15], 混合分布是如下定义的: 假定 X 为一随机变量, 其概率密度函数可以表示为:

$$p(x) = \pi_1 f_1(x) + \pi_2 f_2(x) + \cdots + \pi_k f_k(x)$$

其中 $\pi_j > 0, j = 1, \dots, k$,

$$\pi_1 + \pi_2 + \cdots + \pi_k = 1, \quad f_j(\cdot) \geq 0, \quad \int f_j(x) dx = 1, \quad j = 1, \dots, k,$$

则称 X 是一有限混合分布, $p(\cdot)$ 是有限混合分布的密度函数, 参数 $\pi_1, \pi_2, \dots, \pi_k$ 称为混合权数, $f_1(\cdot), f_2(\cdot), \dots, f_k(\cdot)$ 为混合分布中各个成份(以下称为成分分布)的密度函数, k 为混合分布中所包含的成分分布的个数。

设 θ_j 为分布 $f_j(\cdot)$ 的参数, 那么混合分布的密度函数也可用参数形式表示成:

$$\begin{aligned} p(x; \varphi) &= \pi_1 f_1(x; \theta_1) + \pi_2 f_2(x; \theta_2) + \cdots + \pi_k f_k(x; \theta_k) \\ &= \sum_{j=1}^k \pi_j f_j(x; \theta_j) \end{aligned}$$

其中 $\varphi = (\pi_1, \pi_2, \dots, \pi_k, \theta_1, \theta_2, \dots, \theta_k)$ 为混合分布 X 的参数向量。混合分布的研究内容主要是探讨如何由随机变量 X 经分离分析去获得其中所包含的各个成分分布的特征, 联系本研究目的, 是从分离世代的分布中区别出各主基因型成分分布并估计其数字特征。

2.1 混合分布的似然函数

假定 $x = (x_1, x_2, \dots, x_n)$ 是来自总体 X 的一组样本, 那么样本的似然函数为:

$$L(\varphi) = \prod_{i=1}^n p(x_i; \varphi) = \prod_{i=1}^n \sum_{j=1}^k \pi_j f_j(x_i; \theta_j)$$

极大似然估计一般通过求解对数似然函数的极大值点获得, 在以后的分析中 $L(\varphi)$ 表示对数似然函数, 即:

$$L(\varphi) = \sum_{i=1}^n \log p(x_i; \varphi)$$

2.2 混合分布中成分分布个数的估计

混合分布中所包含成分分布个数 k 是一个最基本的参数, 确定成分分布数目方法可分为图形方法和统计检验方法^[15]两大类, 本文利用 Akaike^[7]提出的最大熵(信息)准则(Akaike's information criterion, AIC), 从不同假设中选择一个最优假设, 从而确定参数 k 。AIC 值定义为:

$$AIC = AIC(k) = -2 L(\hat{\varphi}) + 2 N(k)$$

其中 $\hat{\varphi}$ 是假设“ H_0 : 成分分布个数为 k ”时混合分布中参数的极大似然估计, $N(k)$ 是混合模型

中独立参数的个数。通过比较不同 k 值下的 AIC 值，选择使 AIC 值达到最小的 k 值作为成分分布数的估计。

2.3 混合模型中其它参数的估计

混合模型的参数估计可利用 EM(expectation and maximization) 算法来完成，分两个步骤进行^[4, 6, 10]。

E 步骤：计算完全数据似然函数 $L_c(\varphi)$ 在初始值 $\varphi^{(0)}$ 下的期望值 $Q(\varphi, \varphi^{(0)})$ ，

$$\begin{aligned} Q(\varphi, \varphi^{(0)}) &= E\{L_c(\varphi) | X; \varphi^{(0)}\} \\ &= \sum_{i=1}^n \sum_{j=1}^k \tau_{ij}^{(0)} \cdot [\log \pi_j + f_j(x_i; \theta)] \end{aligned}$$

M 步骤：极大化 $Q(\varphi, \varphi^{(0)})$ ，并用极大值点处的 φ 值代替 $\varphi^{(0)}$ 作为下一轮循环的初始值。 $Q(\varphi, \varphi^{(0)})$ 的极大值点由下式确定：

$$\begin{aligned} \pi_j &= \sum_{i=1}^n \tau_{ij}^{(0)} / n \quad j = 1, \dots, k \\ \sum_{i=1}^n \sum_{j=1}^k \tau_{ij}^{(0)} \partial \log f_j(x_i; \theta) / \partial \theta &= 0 \end{aligned}$$

EM 算法有以下优良性质：(1) M 步骤在正态混合分布的情况下期望函数的极大值点可以用数学式子明确表示出来，而一般情况下企图通过对似然函数求导数来获得极大似然估计的明显数学表示是很困难的；(2) EM 迭代过程中，似然函数是单调增加的，即： $L(\varphi^{(1)}) \geq L(\varphi^k)$ ， $k \geq 0$ 表示第 k 次迭代。这意味着不论对于怎样的初始值，EM 算法最终总能获得一个极大值点。EM 算法的不足之处是收敛速度较慢，并且收敛速度对初始值的选择有一定的依赖性。

对于主基因-多基因混合遗传性状来说，假定成分分布为正态分布，用 μ_j 和 σ_j^2 表示分布的均值和方差，EM 算法具体过程是：

假定 $\varphi^{(0)} = (\pi_1^{(0)}, \dots, \pi_k^{(0)}, \mu_1^{(0)}, \dots, \mu_k^{(0)}, \sigma_1^{2(0)}, \dots, \sigma_k^{2(0)})$

是初始值， $f(x; \mu, \sigma^2)$ 表示正态分布的密度函数，则：

$$\begin{aligned} \tau_{ij}^{(0)} &= Pro(x_i \in G_j | x_i; \varphi^{(0)}) \\ &= \pi_j^{(0)} f(x_i; \mu_j^{(0)}, \sigma_j^{2(0)}) / \sum_{t=1}^k \pi_t^{(0)} f(x_i; \mu_t^{(0)}, \sigma_t^{2(0)}) \\ \pi_j^{(1)} &= \sum_{i=1}^n \tau_{ij}^{(0)} / n \\ \mu_j^{(1)} &= \sum_{i=1}^n \tau_{ij}^{(0)} x_i / (n \pi_j^{(1)}) \\ \sigma_j^{2(1)} &= \sum_{i=1}^n \tau_{ij}^{(0)} (x_i - \mu_j^{(1)})^2 / (n \pi_j^{(1)}) \end{aligned}$$

将求得的参数值 $\varphi^{(1)}$ 代替 $\varphi^{(0)}$ 开始下一轮 EM 迭代。在回交世代中，成分分布还具有相同的方差，用 σ^2 表示，这时，

$$\sigma^{2(1)} = \sum_{i=1}^n \sum_{j=1}^k \tau_{ij}^{(0)} (x_i - \mu_j^{(1)})^2 / n$$

其它参数的表示式同上。

3 利用回交群体鉴定主基因-多基因混合遗传模型

3.1 主基因存在的鉴定及遗传效应的估计

若回交世代均为单一分布，则一般不存在主基因，此处不再讨论。以下分析一个主基因位点(A-a)的情形。假定亲本基因型为AA(P_1)和aa(P_2)， $B_1 = F_1 \times P_1$, $B_2 = F_1 \times P_2$ 。

3.1.1 主基因表现为完全显性 如果通过对回交群体的分解发现回交群体 B_1 是单一正态分布， B_2 是两个正态分布的混合并且两个分布所占比例符合1:1的分离比，此时认为存在有一个主基因(用A-a表示)位点，并表现为完全显性， B_2 群体中两个成分分布的主基因基因型分别为aa和Aa，分布的均值用 μ_1 和 μ_2 表示， μ_1 和 μ_2 与主基因的加性效应 d 和群体平均数 m 的关系为：

$$\begin{aligned}\mu_1 &= m - d \\ \mu_2 &= m + d, \quad \mu_1 < \mu_2\end{aligned}$$

根据极大似然估计的线性不变性，便可获得 m 和 d 的极大似然估计：

$$\begin{aligned}\hat{m} &= 0.5 \hat{\mu}_1 + 0.5 \hat{\mu}_2 \\ \hat{d} &= 0.5 \hat{\mu}_2 - 0.5 \hat{\mu}_1\end{aligned}$$

3.1.2 主基因表现为负向完全显性 如果通过对回交群体的分解发现回交群体 B_2 是单一正态分布， B_1 是两个正态分布的混合并且两个分布所占比例符合1:1的分离比，此时认为只存在有一个主基因(用A-a表示)，并表现为负向完全显性， B_1 群体中两个分布的主基因基因型分别为Aa和AA，分布的均值用 μ_1 和 μ_2 表示， μ_1 和 μ_2 与主基因的加性效应 d 和群体平均数的关系为：

$$\begin{aligned}\mu_1 &= m - d \\ \mu_2 &= m + d\end{aligned}$$

根据极大似然估计的线性不变性，便可获得 m 和 d 的极大似然估计：

$$\begin{aligned}\hat{m} &= 0.5 \hat{\mu}_1 + 0.5 \hat{\mu}_2 \\ \hat{d} &= 0.5 \hat{\mu}_2 - 0.5 \hat{\mu}_1\end{aligned}$$

3.1.3 主基因表现为部分显性或无显性 如果通过对回交群体的分解发现回交群体 B_1 和 B_2 均为两个正态分布的混合并且两个分布所占比例符合1:1的分离比，此时认为只存在有一个主基因(用A-a表示)位点，并表现为部分显性或无显性，分别用 μ_{11} 、 μ_{12} 、 μ_{21} 和 μ_{22} 表示两个回交世代中分布的均值， B_1 群体中两个分布的主基因基因型分别为Aa和AA， B_2 群体中两个分布的主基因基因型分别为aa和Aa，分布的均值与主基因的加性效应 d 、显性效应 h 、群体平均数 m_1 和 m_2 的关系为：

$$\begin{aligned}\mu_{11} &= m_1 + h \\ \mu_{12} &= m_1 + d \\ \mu_{21} &= m_2 - d \\ \mu_{22} &= m_2 + h\end{aligned}$$

根据极大似然估计的线性不变性，便可获得 m_1 、 m_2 、 d 和 h 的极大似然估计：

$$\begin{aligned}\hat{d} &= 0.5(\hat{\mu}_{12} + \hat{\mu}_{22} - \hat{\mu}_{11} - \hat{\mu}_{21}) \\ \hat{h} &= 0.5(\hat{\mu}_{11} + \hat{\mu}_{22} - \hat{\mu}_{12} - \hat{\mu}_{21})\end{aligned}$$

5 讨论

主-多基因性状在分离世代中表现出多峰现象, 具有混合分布的特征, 因此可以利用混合分布理论鉴定主基因的存在并估计遗传效应。但是不同分离世代鉴定主基因的功效是不同的, F_2 世代的功效最低, 回交世代和 F_{2+3} 家系世代的功效较高; 只利用个别分离世代的信息只能对主基因是否存在进行鉴定, 无法对多基因的存在进行鉴定, 多世代的联合分析方法可以对多基因的存在进行鉴定, 并且可以判断多基因是否满足加显性模型, 并进一步估计多基因的遗传效应。

参 考 文 献

- 1 姜长鉴、莫惠栋, 1995, 作物学报, 21(6), 641~648
- 2 马育华编著, 1982, 植物育种的数量遗传学基础, 江苏科学技术出版社, 南京
- 3 莫惠栋, 1993, 作物学报, 19(1), 1~6
- 4 王建康、盖钧镒, 1995, 生物数学学报, 10(4), 87~92
- 5 王建康、盖钧镒(导师), 1996, 数量性状主基因-多基因混合遗传模型的鉴别和遗传参数估计的研究(博士学位论文), 南京农业大学
- 6 王建康、盖钧镒, 1997, 遗传学报, 24(5), 432~440
- 7 Akaike, H., 1977, In: P. R. Krishnaiah(ed.) Application of Statistics, pp. 27~41, North-Holland Publishing Company, Amsterdam
- 8 Elkind, Y. and A. Cahaner, 1986, Theor. Appl. Genet., 72, 377~383
- 9 Elkind, Y. and A. Cahaner, 1990, Heredity, 64, 205~213
- 10 Dempster, A. P., N. M. Laird and D. B. Robin, 1977, J. R. Statist. Soc. B., 39, 1~38
- 11 Elston, R. C., 1984, Genetics, 108, 733~744
- 12 Goffinet, B., P. Loisel and B. Laurent, 1990, Biometrics, 46, 583~594
- 13 Loisel, P., B. Goffinet, H. Monod and G. M. De Oca, 1994, Biometrics, 50, 512~516
- 14 Mather, K. and J. L. Jinks, 1982, Biometrical Genetics, Chapman and Hall
- 15 McLachlan, G. J., 1988, Mixture Models: Inference and Applications to Clustering, Marcel Dekker, Inc.
- 16 Tanksley, S. D., 1993, Annu. Rev. Genet., 27, 205~233

Identification of Major Gene and Polygene Mixed Inheritance Model from Backcrosses or F_{2+3} Families

Gai Junyi Wang Jiankang

(Soybean Research Institute, Nanjing Agricultural University, Nanjing 210095)

Abstract For major-polygene traits, the distribution of a segregating population demonstrates multimodality which is the characteristic of a mixture of more than one distributions. A statistical method to identify major gene and polygene mixed inheritance from backcrosses and F_{2+3} families was developed by using the theory of mixture distribution. With the method, the existence of major gene and its genetic effects could be determined.

Key words Quantitative trait; Major gene and polygene mixed inheritance; Mixture distribution; EM algorithm