

# METALLICITY

OCTOBER 26, 2010

## 1. MIXTURE MODEL

Given  $n$  observed  $(\frac{\alpha}{\text{Fe}}, \frac{\text{Fe}}{\text{H}})$  metallicities as  $\{(x_i, y_i)\}_{i=1}^n$ , or as  $(\mathbf{x}, \mathbf{y})$ , each of which is drawn from one of  $m$  known model densities. We model the density of observations using the mixture model

$$(1) \quad f(x, y) = \sum_{j=1}^m \pi_j f_j(x, y)$$

where

$$\sum_{j=1}^m \pi_j = 1 \quad \pi_j \geq 0, \quad j = 1, \dots, m$$

From the summation constraint,  $\boldsymbol{\pi}$  has  $m - 1$  free parameters:

$$\boldsymbol{\pi} = (\pi_1, \dots, \pi_{m-1}, 1 - \pi_1 - \dots - \pi_{m-1})$$

Thus the likelihood of (1) is

$$\begin{aligned} L(\boldsymbol{\pi}) &= \prod_{i=1}^n f(x_i, y_i) \\ &= \prod_{i=1}^n \left\{ \sum_{j=1}^m \pi_j f_j(x_i, y_i) \right\} \\ \log L(\boldsymbol{\pi}) &= \sum_{i=1}^n \log \left( \sum_{j=1}^m \pi_j f_j(x_i, y_i) \right) \end{aligned}$$

Maximizing  $\log L(\boldsymbol{\pi})$  with respect to  $\boldsymbol{\pi}$  will yield  $\hat{\boldsymbol{\pi}}_{\text{MLE}}$ , but this arduous task can be avoided by adding a latent indicator,  $z$ , to the observed data  $(\mathbf{x}, \mathbf{y})$ , representing the model group from which that observation was generated. Let  $G_j$  be the  $j^{\text{th}}$  model group, and let

$$z_{ij} = \mathbf{1}\{(x_i, y_i) \mapsto G_j\}$$

The complete data likelihood is defined over the complete data  $\{(x_i, y_i, \mathbf{z}_i)\}_{i=1}^n$  as

$$L(\boldsymbol{\pi}) = \prod_{i=1}^n \prod_{j=1}^m \left\{ f_j(x_i, y_i) \right\}^{z_{ij}} \pi_j^{z_{ij}}$$

$$(2) \quad l(\boldsymbol{\pi}) = \sum_{i=1}^n \sum_{j=1}^m z_{ij} \log \{ \pi_j f_j(x_i, y_i) \}$$

## 2. EXPECTATION MAXIMIZATION

One way to estimate  $\boldsymbol{\pi}$  is to use a maximum likelihood estimate,  $\hat{\boldsymbol{\pi}}$ , computed using expectation maximization. Starting from an initial set of guesses,  $\boldsymbol{\pi}^{(0)}$ , we iteratively find the expected value of the likelihood, (2), conditional on the data, and then find the  $\text{argmax}_{\boldsymbol{\pi}}$  of this expectation. The maximizing value the  $t^{\text{th}}$  iteration,  $\hat{\boldsymbol{\pi}}^{(t)}$ , is then used as the starting value for the next run, and we continue until the likelihood changes by less than  $10^{-3}$  over twenty five iterations.

**2.1. Expectation step.** First we find the expected value of the log likelihood, (2), conditional on the data. Note that since  $z_{ij}$  is an indicator function, its expected value is equal to the probability that data point  $i$  comes from model  $j$ .

$$\mathbb{E}_{\boldsymbol{\pi}}[l(\boldsymbol{\pi})|\mathbf{x}, \mathbf{y}] = \sum_{i=1}^n \sum_{j=1}^m \mathbb{E}_{\boldsymbol{\pi}}[z_{ij}|x_i, y_i] \{ \log f_j(x_i, y_i) + \log \pi_j \}$$

Since we're ultimately maximizing, the non-constant component is of primary interest, and can be analytically specified by applying Bayes' rule:

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\pi}}[z_{ij}|x_i, y_i] &= \text{Probability}\left((x_i, y_i) \mapsto G_j | x_i, y_i\right) \\ &= \Pr_{\boldsymbol{\pi}}(z_{ij}|x_i, y_i) \\ &= \frac{p(x_i, y_i | z_{ij} = 1) p(z_{ij} = 1)}{p(x_i, y_i)} \end{aligned}$$

Thus the expected value of the indicator variable,  $z_{ij}$ , given the data and the parameters,  $\boldsymbol{\pi}$ , of the data's distribution defined by (1) is

$$(3) \quad \mathbb{E}_{\boldsymbol{\pi}}[z_{ij}|x_i, y_i] = \frac{\pi_j f_j(x_i, y_i)}{\sum_{j=1}^m \pi_j f_j(x_i, y_i)}$$

To iteratively evaluate this expectation, we let  $w_{ij}^{(t)}$  be (3) at the  $t^{\text{th}}$  step:

$$w_{ij}^{(t+1)} = \begin{cases} \frac{\pi_j^{(t)} f_j(x_i, y_i)}{\sum_{k=1}^m \pi_k^{(t)} f_k(x_i, y_i)} & j = 1, \dots, m-1 \\ 1 - w_{i1} - \dots - w_{i,m-1} & j = m \end{cases}$$

Since  $\boldsymbol{\pi}$  is not defined for the first evaluation, we use a random initialization to generate  $\mathbf{w}_j^{(0)}$ . Convergence is not sensitive to the choice of values in this case, but may be if the likelihood is riddled with local maxima.

## 2.2. Maximizing with respect to $\boldsymbol{\pi}$ .

$$\begin{aligned} 0 &= \frac{\partial}{\partial \pi_k} \mathbb{E}_{\boldsymbol{\pi}}[l(\boldsymbol{\pi})|\mathbf{x}, \mathbf{y}] \\ &= \sum_{i=1}^n \left\{ w_{ij}^{(0)} \frac{1}{\pi_k} - w_{im}^{(0)} \frac{1}{1 - \pi_1 - \dots - \pi_{m-1}} \right\}, k = 1, \dots, m-1 \end{aligned}$$

as  $\pi_m = 1 - \pi_1 - \dots - \pi_{m-1}$ .

$$\begin{aligned} \frac{1}{\pi_1} \sum_{i=1}^n w_{i1}^{(0)} &= \dots = \frac{1}{\pi_{m-1}} \sum_{i=1}^n w_{i,m-1}^{(0)} = c \\ \hat{\pi}_k &= \frac{\sum_{i=1}^n w_{ik}^{(0)}}{c} \\ \pi_j^{(1)} &= \frac{\sum_{i=1}^n w_{ij}^{(0)}}{n} \end{aligned}$$

And in general,

$$\pi_j^{(t+1)} = \frac{\sum_{i=1}^n w_{ij}^{(t)}}{n}$$

## 3. COVARIANCE

The asymptotic covariance matrix of  $\hat{\boldsymbol{\pi}}$  can be approximated by the inverse of the observed Fisher information matrix. Since there are only  $m - 1$  free parameters, let  $\boldsymbol{\pi}' = (\pi_1, \dots, \pi_{m-1})$ . The likelihood can then be expressed as:

$$(4) \quad l(\boldsymbol{\pi}') = \sum_{i=1}^n \log \left\{ \left( \sum_{j=1}^{m-1} \pi_j f_j \right) + (1 - \pi_1, \dots, \pi_{m-1}) f_m \right\}$$

The observed information matrix,  $I(\boldsymbol{\pi}'|\mathbf{x}, \mathbf{y})$ , is given by the  $m - 1 \times m - 1$  negative hessian of (4):

$$I(\boldsymbol{\pi}'|\mathbf{x}, \mathbf{y}) = -\frac{\partial^2 l(\boldsymbol{\pi}')}{\partial \boldsymbol{\pi}' \partial \boldsymbol{\pi}'^T} = - \begin{bmatrix} \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial^2 \pi_1} & \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial \pi_1 \partial \pi_2} & \cdots & \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial \pi_1 \partial \pi_{m-1}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial \pi_{m-1} \partial \pi_1} & \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial \pi_{m-1} \partial \pi_2} & \cdots & \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial^2 \pi_{m-1}} \end{bmatrix}$$

where

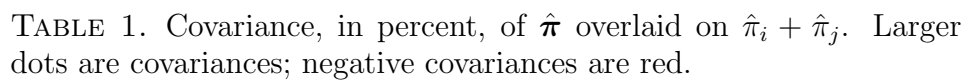
$$\begin{aligned} \frac{\partial l(\boldsymbol{\pi}')}{\partial \pi_k} &= \sum_{i=1}^n \frac{f_k - f_m}{\sum_{j=1}^m \pi_j f_j} \\ \frac{\partial^2 l(\boldsymbol{\pi}')}{\partial \pi_k \partial \pi_r} &= - \sum_{i=1}^n \frac{(f_k - f_m)(f_r - f_m)}{(\sum_{j=1}^m \pi_j f_j)^2 f_r} \end{aligned}$$

The inverse of  $I(\boldsymbol{\pi}'|\mathbf{x}, \mathbf{y})$  provides estimates of the variance, covariance, and correlation of  $\hat{\boldsymbol{\pi}}$  as

$$\text{Cov}(\hat{\pi}_p, \hat{\pi}_q) = \begin{cases} [I^{-1}(\hat{\boldsymbol{\pi}})]_{pq} & p, q < m \\ - \sum_{j=1}^{m-1} \text{Cov}(\hat{\pi}_j, \hat{\pi}_q) & p = m, q < m \\ \sum_{j=1}^{m-1} \sum_{k=1}^{m-1} \text{Cov}(\hat{\pi}_j, \hat{\pi}_q) & p, q = m \end{cases}$$

$$\text{Var}(\hat{\pi}_j) = \sigma_j^2 = \left\{ \text{Cov}(\hat{\boldsymbol{\pi}}) \right\}_{jj}$$

$$\text{Corr}(\hat{\pi}_p, \hat{\pi}_q) = \frac{\text{Cov}(\hat{\pi}_p, \hat{\pi}_q)}{\sqrt{\sigma_p^2 \sigma_q^2}}$$



0.0001	0.001	-	-	-0.0004	-0.0001	-	-	0.0001	-	-	-	-	-	-	0.0007
-0.001	-0.0009	-	-0.0002	0.0015	-0.0004	-	-	-	-	-	-	-	-	-	-0.001
-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
-0.0002	0.0003	-	-	-	-	-	-	-	-	-	-	-	-	-	-
-	-0.0001	-	-	0.0001	-	-	-	-	-	-	-	-	-	-	-
0.0003	-0.0001	-	0.0001	-0.0004	0.0002	-	-	-	-	-	-	-	-	-	-
-0.0001	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
-0.0011	-0.0001	-	-0.0013	0.0053	-0.0021	-	-	-	-	-	-	0.0001	-	-	0.001
0.0023	0.0013	-0.0001	0.0008	-0.002	-	-0.0001	0.0001	-	-	-	-	-	-	-	0.0023
-	-	-	-	-0.0001	0.0001	-	-	-	-	-	-	-	-	-	-
-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
0.0003	-0.0003	-	-	-0.0001	0.0001	-	-	-	-	-	-	-	-	-	-
-0.0058	0.0041	0.0001	-0.0019	0.003	-0.0011	0.0002	0.0001	0.0003	-	-	-0.0001	-0.0001	-0.0001	-	-0.0012
0.0076	-0.0053	-0.0001	0.0019	-0.0086	0.0045	-0.0003	0.0002	-0.0003	-	-	-	-	-0.0001	-	-0.0005
0.0025	-0.0001	-0.0001	-0.0005	-0.0018	0.0011	-0.0001	0.0003	-	-	-	-	-	-	-	0.0012

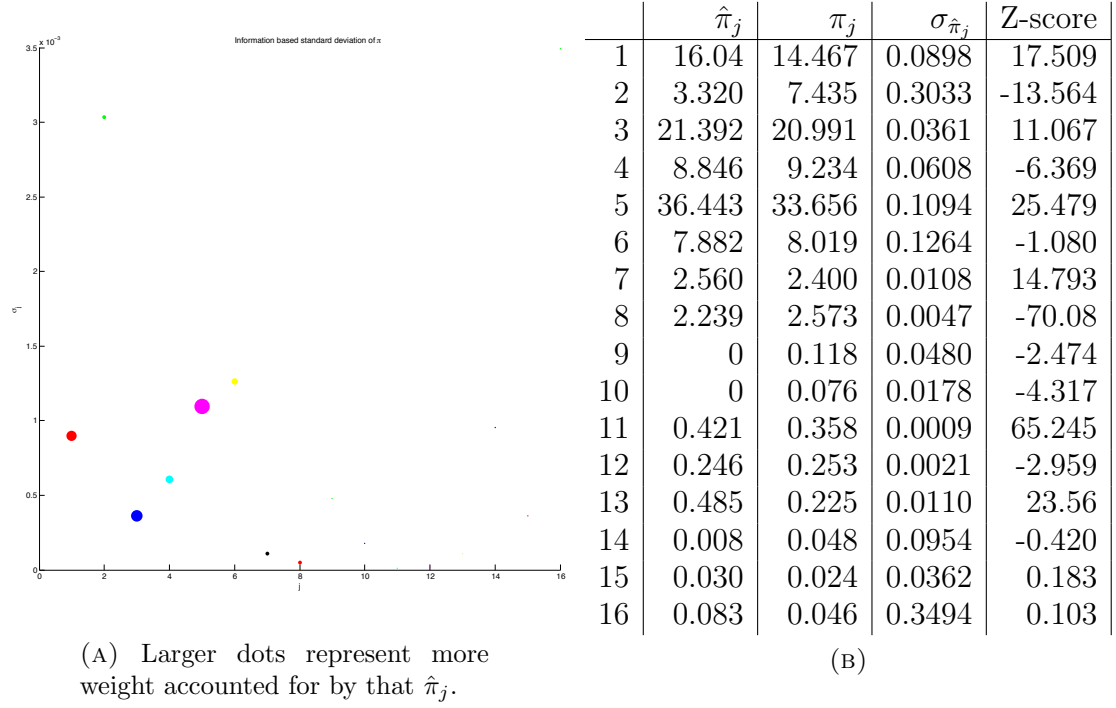


FIGURE 1. Variance of  $\hat{\pi}$ , true  $\pi$ , in percent, and z-scores.

$$\text{Z-score} = \frac{\hat{\pi}_j - \pi}{\sigma_{\hat{\pi}_j}}$$

## 4. LIKELIHOOD RATIO TEST

Given certain conditions

$$H_0 : \boldsymbol{\pi} = \boldsymbol{\pi}_{\text{true}}$$

$$H_1 : \boldsymbol{\pi} \neq \boldsymbol{\pi}_{\text{true}}$$

$$\Lambda = -2 \log \frac{\sup_{\boldsymbol{\pi}=\boldsymbol{\pi}_{\text{true}}} l(\boldsymbol{\pi})}{\sup_{\boldsymbol{\pi}} l(\boldsymbol{\pi})} = -2 \{l(\boldsymbol{\pi}_{\text{true}}) - l(\hat{\boldsymbol{\pi}})\} \sim \chi_{m-1}^2$$

$$\Lambda_{\text{Halo } 3} = 25.025 \sim \chi_{15}^2$$

$$\text{p-value } 10\text{k} = 4.961\%$$

$$\text{p-value } 30\text{k} = 1.3 \times 10^{-5}\%$$

$$\text{p-value } 50\text{k} = 1.4 \times 10^{-12}\%$$

Thus we accept  $H_0$  when requiring 95% or less confidence; there is only a 4.961% chance we would see a value this extreme or more given  $H_0$  is true. This holds for 400 to 1600 EM iterations.