

# Introduction to Bayesian Statistics - 5

*Edoardo Milotti*

Università di Trieste and INFN-Sezione di Trieste

# Two important computational techniques with a Bayesian basis

1. The EM algorithm
2. Image processing techniques (MLM, MEM, etc.)

# 1. The EM algorithm (Dempster, Laird & Rubin, 1977)

Recall the max. likelihood principle:

$$\begin{aligned} P(\boldsymbol{\theta} \mid \mathbf{d}, I) &= \frac{P(\mathbf{d} \mid \boldsymbol{\theta}, I)}{P(\mathbf{d} \mid I)} \cdot P(\boldsymbol{\theta} \mid I) \\ &= \frac{\mathcal{L}(\mathbf{d}, \boldsymbol{\theta})}{P(\mathbf{d} \mid I)} \cdot P(\boldsymbol{\theta} \mid I) \propto \mathcal{L}(\mathbf{d}, \boldsymbol{\theta}) \end{aligned}$$

uniform distribution  
(usually an improper prior)

evidence

likelihood

in this (approximate) setting, the MAP estimate coincides with the ML estimate.

when data are independent and identically distributed (i.i.d.)  
we find the following likelihood function

$$\mathcal{L}(\mathbf{d}, \boldsymbol{\theta}) = \prod_i p(d_i | \boldsymbol{\theta})$$

and we estimate the parameters by maximizing the likelihood function

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{d}, \boldsymbol{\theta})$$

or, equivalently, its logarithm

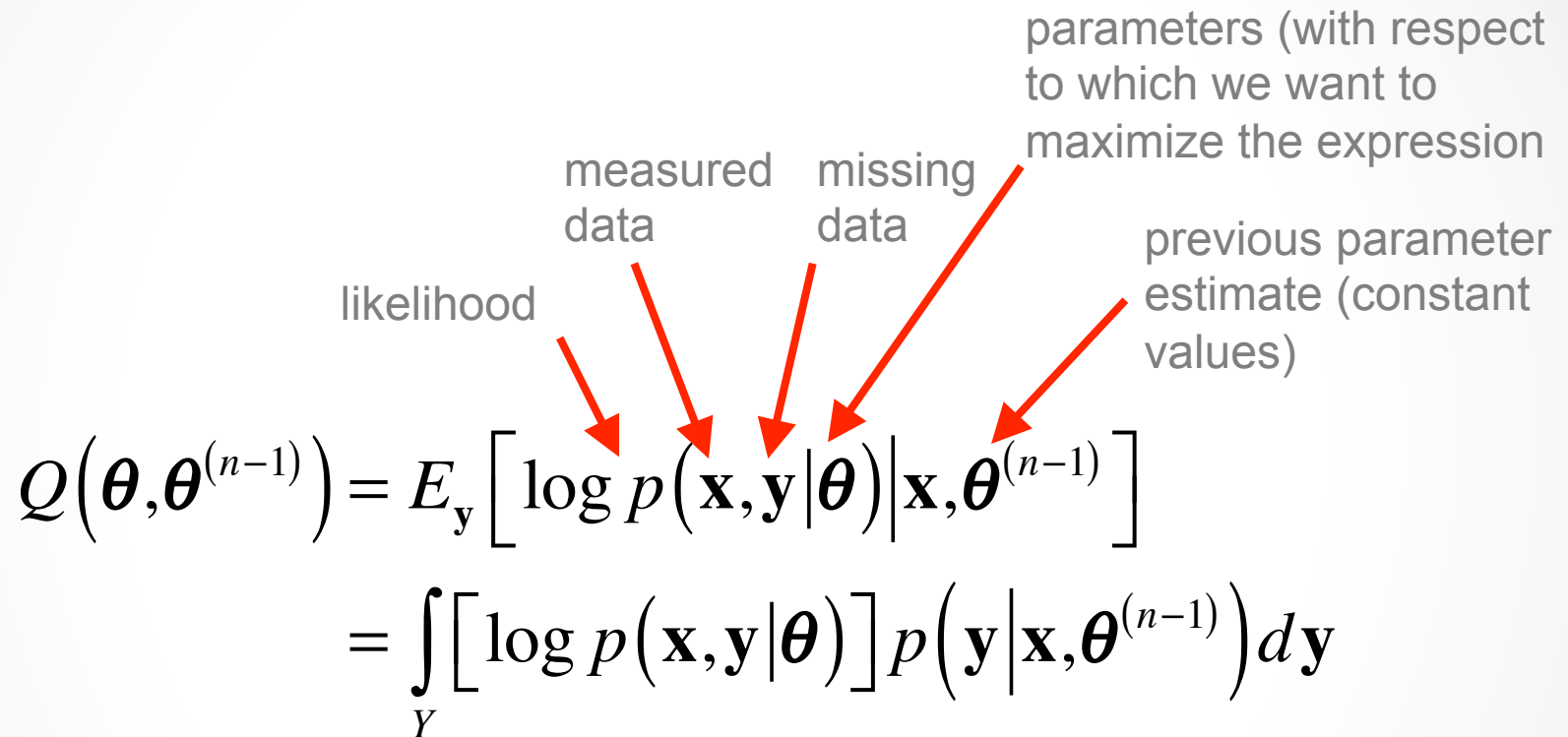
$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} [\log \mathcal{L}(\mathbf{d}, \boldsymbol{\theta})]$$

(in real life, this procedure is often complex and almost invariably it requires a numerical solution)



*The EM algorithm is used to maximize likelihood with incomplete information, and it has two main steps that are iterated until convergence:*

**E. expectation of the log-likelihood, averaged with respect to missing data:**


$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(n-1)}) = E_{\mathbf{y}} \left[ \log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \middle| \mathbf{x}, \boldsymbol{\theta}^{(n-1)} \right]$$
$$= \int_Y \left[ \log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \right] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}^{(n-1)}) d\mathbf{y}$$

**M. maximization of the averaged log-likelihood with respect to parameters:**

$$\boldsymbol{\theta}^{(n)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(n-1)})$$

## Example: an experiment with an exponential model (Flury and Zoppè)

Light bulbs fail following an exponential distribution with mean failure time  $\theta$

*To estimate the mean two experiments are performed*

1.  $n$  light bulbs are tested, all failure times  $u_i$  are recorded
2.  $m$  light bulbs are tested, only the total number  $r$  of bulbs failed at time  $t$  are recorded

$$1. \quad \mathcal{L} = \prod_{i=1}^n \frac{1}{\theta} \exp\left(-\frac{u_i}{\theta}\right) = \frac{1}{\theta^n} \exp\left(-\frac{\sum_i u_i}{\theta}\right) = \frac{1}{\theta^n} \exp\left(-\frac{n\langle u \rangle}{\theta}\right)$$

$$2. \quad \mathcal{L} = \prod_{i=1}^m \frac{1}{\theta} \exp\left(-\frac{v_i}{\theta}\right)$$

missing data!

combined likelihood

$$\frac{1}{\theta^n} \exp\left(-\frac{n\langle u \rangle}{\theta}\right) \cdot \prod_{i=1}^m \frac{1}{\theta} \exp\left(-\frac{v_i}{\theta}\right)$$

log-likelihood

$$-n \ln \theta - \frac{n\langle u \rangle}{\theta} - \sum_{i=1}^m \left( \ln \theta + \frac{v_i}{\theta} \right)$$

expected failure time for a bulb  
that is still burning at time  $t$

$$t + \theta$$

expected failure time for a bulb  
that is not burning at time  $t$

$$\theta - \frac{t \exp(-t/\theta)}{1 - \exp(-t/\theta)}$$

# Note on mean failure time for a bulb that is not burning at time $t$

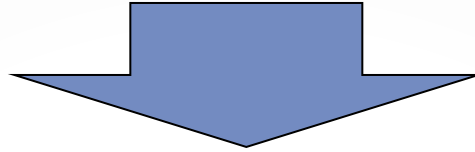
$$p(t') \propto \frac{1}{\theta} e^{-t'/\theta} \quad 0 \leq t' \leq t$$

$$\text{normalization} = \int_0^t p(t') dt' = \int_0^t \frac{dt'}{\theta} e^{-t'/\theta} = 1 - e^{-t/\theta}$$

$$\text{mean failure time} = \int_0^t t' p(t') dt' = \frac{1}{1 - e^{-t/\theta}} \int_0^t t' e^{-t'/\theta} \frac{dt'}{\theta}$$

$$= \frac{\theta}{1 - e^{-t/\theta}} \left[ 1 - e^{-t/\theta} - (t/\theta) e^{-t/\theta} \right]$$

$$= \theta - \frac{te^{-t/\theta}}{1 - e^{-t/\theta}}$$



average log-likelihood

$$\begin{aligned} Q &= E \left[ -n \ln \theta - \frac{n \langle u \rangle}{\theta} + \sum_{i=1}^m \left( -\ln \theta - \frac{v_i}{\theta} \right) \right] \\ &= -(n+m) \ln \theta - \frac{n \langle u \rangle}{\theta} - \frac{r}{\theta} \left( \theta - \frac{t \exp(-t/\theta)}{1 - \exp(-t/\theta)} \right) - \frac{(m-r)}{\theta} (\theta + t) \end{aligned}$$

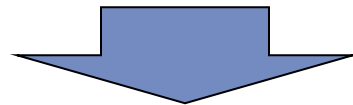
*this ends the expectation step*

the max of the mean likelihood

$$Q = -(n+m)\ln\theta - \frac{1}{\theta} \left[ n\langle u \rangle + r \left( \theta - \frac{t \exp(-t/\theta)}{1 - \exp(-t/\theta)} \right) + (m-r)(\theta + t) \right]$$

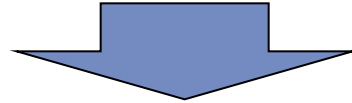
can be found by maximizing the approximate expression

$$Q \approx -(n+m)\ln\theta - \frac{1}{\theta} \left[ n\langle u \rangle + r \left( \theta^{(k)} - \frac{t \exp(-t/\theta^{(k)})}{1 - \exp(-t/\theta^{(k)})} \right) + (m-r)(\theta^{(k)} + t) \right]$$



$$\frac{dQ}{d\theta} \approx -(n+m)\frac{1}{\theta} + \frac{1}{\theta^2} \left[ n\langle u \rangle + r \left( \theta^{(k)} - \frac{t \exp(-t/\theta^{(k)})}{1 - \exp(-t/\theta^{(k)})} \right) + (m-r)(\theta^{(k)} + t) \right] = 0$$

$$\frac{dQ}{d\theta} \approx -(n+m)\frac{1}{\theta} + \frac{1}{\theta^2} \left[ n\langle u \rangle + r \left( \theta^{(k)} - \frac{t \exp(-t/\theta^{(k)})}{1 - \exp(-t/\theta^{(k)})} \right) + (m-r)(\theta^{(k)} + t) \right] = 0$$

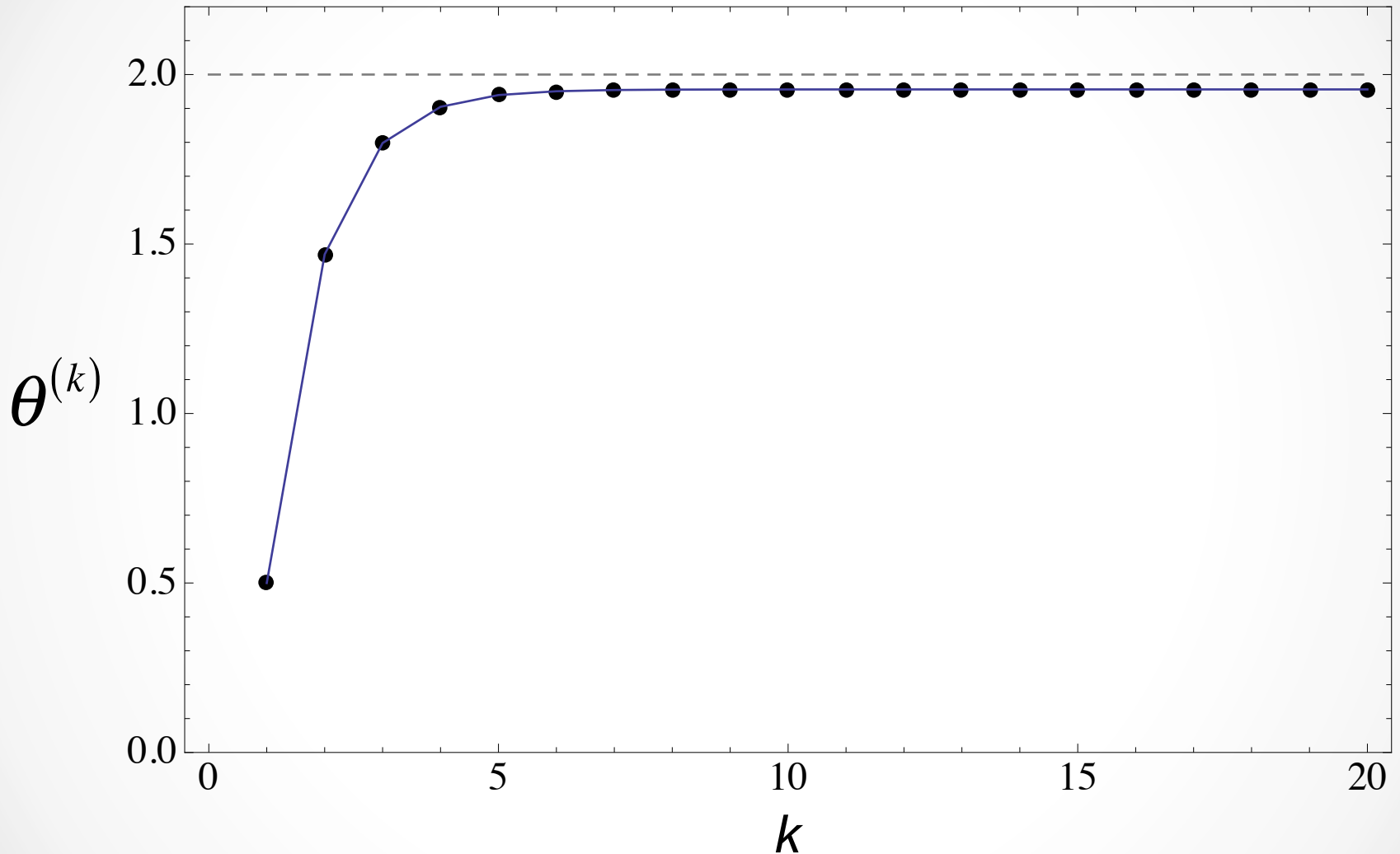


$$\theta^{(k+1)} = \frac{1}{n+m} \left[ n\langle u \rangle + r \left( \theta^{(k)} - \frac{t \exp(-t/\theta^{(k)})}{1 - \exp(-t/\theta^{(k)})} \right) + (m-r)(\theta^{(k)} + t) \right]$$

iterate this until convergence ...



Example with mean failure time = 2 (a.u.), and randomly generated data ( $n = 100$ ;  $m = 100$ ). In this example  $r = 36$ .

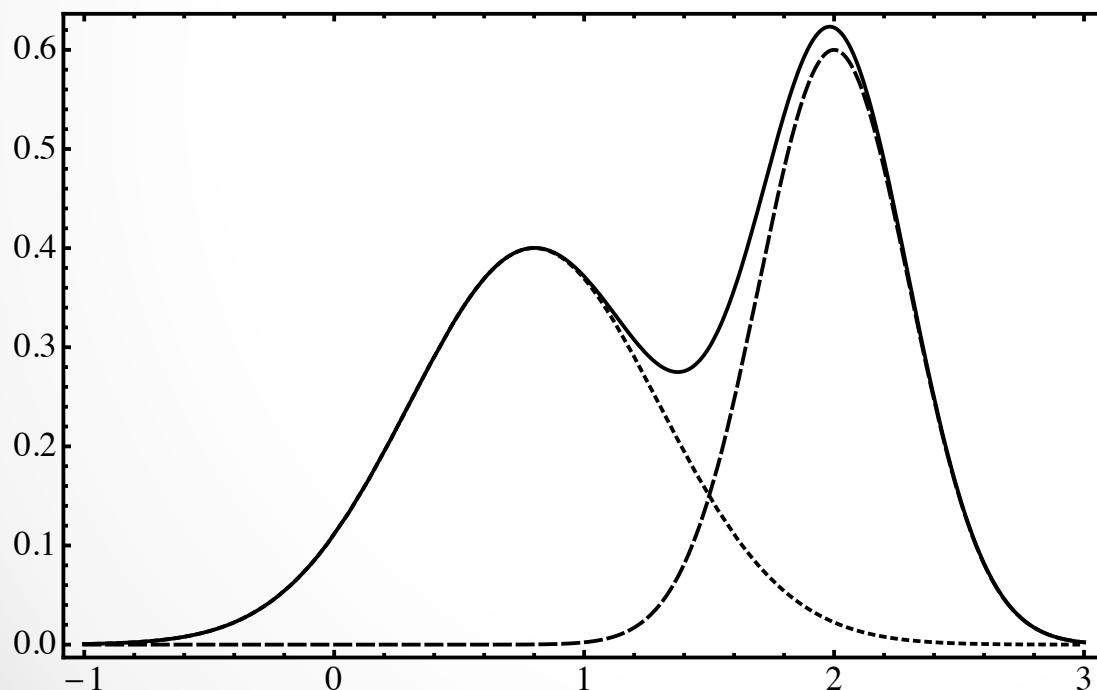


Important application of the EM method: parameters of “mixture models”.

$$p(x_n|\boldsymbol{\theta}) = \sum_{i=1}^M \alpha_i p_i(x_n|\boldsymbol{\theta}_i)$$

$$\boldsymbol{\theta} = (\alpha_1, \dots, \alpha_M; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M)$$

$$\sum_{i=1}^M \alpha_i = 1$$



Example: a Gaussian mixture model (M=2)

direct maximization of log likelihood

$$\begin{aligned}\log \mathcal{L}(\mathbf{x}, \boldsymbol{\theta}) &= \log \prod_n p(x_n | \boldsymbol{\theta}) = \sum_n \log p(x_n | \boldsymbol{\theta}) \\ &= \sum_n \log \left[ \sum_{i=1}^M \alpha_i p_i(x_n | \boldsymbol{\theta}_i) \right]\end{aligned}$$

this is difficult ... however we can do it differently with a reinterpretation of the mixture model parameters ...

$\alpha_k$  = probability of drawing the  $k$ -th component of the mixture model



new (hidden) variable:  $y = k$  = index of component

thus we must redefine data and parameters

new likelihood which includes the hidden variables

$$\begin{aligned}\log \mathcal{L}'(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) &= \log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \\ &= \log \prod_n p(x_n, y_n | \boldsymbol{\theta}) \\ &= \sum_n \log \left[ p(x_n | y_n, \boldsymbol{\theta}) p(y_n | \boldsymbol{\theta}) \right] \\ &= \sum_n \log \left[ \alpha_{y_n} p_{y_n} \left( x_n | \boldsymbol{\theta}_{y_n} \right) \right]\end{aligned}$$

(  $\boldsymbol{\theta}_i$  are the parameters restricted to the i-th component)

The structure is simpler now, however there are hidden variables.

Now we proceed by averaging the likelihood  
(Expectation step)

new parameter  
estimate

previous parameter  
estimate

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) &= E_{\mathbf{y}} \left[ \log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) | \mathbf{x}, \boldsymbol{\theta}^{(i-1)} \right] \\ &= \int_Y \left[ \log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \right] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}^{(i-1)}) d\mathbf{y} \\ &\rightarrow \sum_{\mathbf{y}} \left[ \log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \right] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}^{(i-1)}) \end{aligned}$$

sum instead of integral, because the  
 $\mathbf{y}$  variates are discrete

prior probabilities in the expression of the averaged log-likelihood

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) = \sum_{\mathbf{y}} [\log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta})] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}^{(i-1)})$$

where:

$$p(y_n | x_n, \boldsymbol{\theta}) = \frac{p(x_n | y_n, \boldsymbol{\theta}) p(y_n | \boldsymbol{\theta})}{p(x_n | \boldsymbol{\theta})} = \frac{\alpha_{y_n} p_{y_n}(x_n | \boldsymbol{\theta}_{y_n})}{\sum_{k=1}^M \alpha_k p_k(x_n | \boldsymbol{\theta}_k)}$$

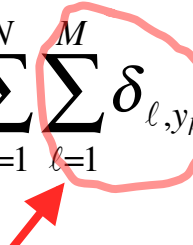
$$p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) = \prod_{n=1}^N p(y_n | x_n, \boldsymbol{\theta}) = \prod_{n=1}^N \frac{\alpha_{y_n} p_{y_n}(x_n | \boldsymbol{\theta}_{y_n})}{\sum_{k=1}^M \alpha_k p_k(x_n | \boldsymbol{\theta}_k)}$$

Therefore, using  $\log \mathcal{L}'(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \sum_n \log \left[ \alpha_{y_n} p_{y_n} \left( x_n \mid \boldsymbol{\theta}_{y_n} \right) \right]$

$$\text{and } p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) = \prod_{n=1}^N p(y_n \mid x_n, \boldsymbol{\theta})$$

we find

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) &= \sum_{\mathbf{y}} [\log p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\theta})] p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}^{(i-1)}) \\ &= \sum_{\mathbf{y}} \sum_{k=1}^N \log \left[ \alpha_{y_k} p_{y_k} \left( x_k \mid \boldsymbol{\theta}_{y_k} \right) \right] \prod_{j=1}^N p(y_j \mid x_j, \boldsymbol{\theta}^{(i-1)}) \\ &= \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{k=1}^N \log \left[ \alpha_{y_k} p_{y_k} \left( x_k \mid \boldsymbol{\theta}_{y_k} \right) \right] \prod_{j=1}^N p(y_j \mid x_j, \boldsymbol{\theta}^{(i-1)}) \end{aligned}$$

$$\begin{aligned}
 Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) &= \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{k=1}^N \log \left[ \alpha_{y_k} p_{y_k} \left( x_k \middle| \boldsymbol{\theta}_{y_k} \right) \right] \prod_{j=1}^N p \left( y_j \middle| x_j, \boldsymbol{\theta}^{(i-1)} \right) \\
 &= \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{k=1}^N \sum_{\ell=1}^M \delta_{\ell, y_k} \log \left[ \alpha_{\ell} p_{\ell} \left( x_k \middle| \boldsymbol{\theta}_{\ell} \right) \right] \prod_{j=1}^N p \left( y_j \middle| x_j, \boldsymbol{\theta}^{(i-1)} \right)
 \end{aligned}$$


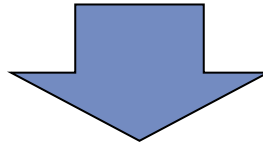
to decouple the variables, we add one sum and one Kronecker's delta...

after the decoupling, we can use the normalization of conditional probabilities

$$\sum_{y_j=1}^M p \left( y_j \middle| x_j, \boldsymbol{\theta}^{(i-1)} \right) = 1$$



$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) = \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{k=1}^N \sum_{\ell=1}^M \delta_{\ell, y_k} \log[\alpha_{\ell} p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] \prod_{j=1}^N p(y_j | x_j, \boldsymbol{\theta}^{(i-1)})$$



$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) &= \sum_{\ell=1}^M \sum_{k=1}^N \log[\alpha_{\ell} p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \delta_{\ell, y_k} \prod_{n=1}^N p(y_n | x_n, \boldsymbol{\theta}^{(i-1)}) \\ &= \sum_{\ell=1}^M \sum_{k=1}^N \log[\alpha_{\ell} p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] \left\{ \sum_{y_1=1}^M \dots \sum_{y_{k-1}=1}^M \sum_{y_{k+1}=1}^M \dots \sum_{y_N=1}^M \prod_{\substack{j=1 \\ j \neq k}}^N p(y_j | x_j, \boldsymbol{\theta}^{(i-1)}) \right\} p(\ell | x_k, \boldsymbol{\theta}^{(i-1)}) \\ &= \sum_{\ell=1}^M \sum_{k=1}^N \log[\alpha_{\ell} p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] \left\{ \prod_{\substack{j=1 \\ j \neq k}}^N \sum_{y_j=1}^M p(y_j | x_j, \boldsymbol{\theta}^{(i-1)}) \right\} p(\ell | x_k, \boldsymbol{\theta}^{(i-1)}) \\ &= \sum_{\ell=1}^M \sum_{k=1}^N \log[\alpha_{\ell} p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] p(\ell | x_k, \boldsymbol{\theta}^{(i-1)}) \end{aligned}$$



these sums all add to 1  
(normalization of conditional probabilities)

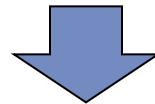
$$\begin{aligned}
 Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i-1)}) &= \sum_{\ell=1}^M \sum_{k=1}^N \log [\alpha_{\ell} p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] p(\ell | x_k, \boldsymbol{\theta}^{(i-1)}) \\
 &= \sum_{\ell=1}^M \sum_{k=1}^N [\log \alpha_{\ell}] p(\ell | x_k, \boldsymbol{\theta}^{(i-1)}) + \sum_{\ell=1}^M \sum_{k=1}^N [\log p_{\ell}(x_k | \boldsymbol{\theta}_{\ell})] p(\ell | x_k, \boldsymbol{\theta}^{(i-1)})
 \end{aligned}$$

this depends only on the  $\alpha$  parameters

this term depends on the parameters of the component distributions

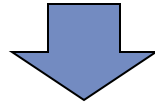
Thus there are two terms that can be maximized separately.  
Moreover, the first term must be maximized with the normalization constraint, i.e.

$$\frac{\partial}{\partial \alpha_m} \left[ \sum_{\ell=1}^M \sum_{k=1}^N \log \alpha_{\ell} p(\ell | x_k, \boldsymbol{\theta}^{(i-1)}) + \lambda \left( \sum_{\ell=1}^M \alpha_{\ell} - 1 \right) \right] = 0$$

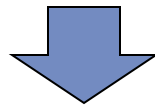


$$\sum_{k=1}^N \frac{1}{\alpha_m} p(m | x_k, \boldsymbol{\theta}^{(i-1)}) + \lambda = 0$$

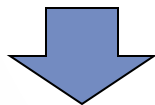
$$\sum_{k=1}^N \frac{1}{\alpha_m} p(m|x_k, \boldsymbol{\theta}^{(i-1)}) + \lambda = 0$$



$$\sum_{k=1}^N p(m|x_k, \boldsymbol{\theta}^{(i-1)}) = -\lambda \alpha_m$$



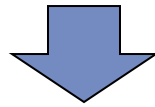
$$\sum_{m=1}^M \sum_{k=1}^N p(m|x_k, \boldsymbol{\theta}^{(i-1)}) = -\lambda \sum_{m=1}^M \alpha_m$$



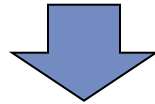
$$\lambda = -N \quad \Rightarrow \quad \alpha_m = \frac{1}{N} \sum_{k=1}^N p(m|x_k, \boldsymbol{\theta}^{(i-1)})$$

This is as far as we can go without introducing an explicit form for the component distributions: now we explicitly consider the 1D Gaussian mixture model:

$$p_{\ell}(x|\mu_{\ell},\sigma_{\ell}) = \frac{1}{\sqrt{2\pi\sigma_{\ell}^2}} \exp\left(-\frac{(x-\mu_{\ell})^2}{2\sigma_{\ell}^2}\right)$$

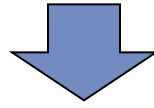


$$\sum_{\ell=1}^M \sum_{k=1}^N \log p_{\ell}(x_k|\boldsymbol{\theta}_{\ell}) p(\ell|x_k, \boldsymbol{\theta}^{(i-1)}) = \sum_{\ell=1}^M \sum_{k=1}^N \left[ -\frac{1}{2} \ln(2\pi\sigma_{\ell}^2) - \frac{(x_k - \mu_{\ell})^2}{2\sigma_{\ell}^2} \right] p(\ell|x_k, \mu_{\ell}^{(i-1)}, \sigma_{\ell}^{(i-1)})$$



$$\frac{\partial}{\partial \mu_m} \sum_{\ell=1}^M \sum_{k=1}^N \log p_{\ell}(x_k|\boldsymbol{\theta}_{\ell}) p(\ell|x_k, \boldsymbol{\theta}^{(i-1)}) = -2 \sum_{k=1}^N \frac{(x_k - \mu_m)}{2\sigma_m^2} p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)}) = 0$$

$$\frac{\partial}{\partial \mu_m} \sum_{\ell=1}^M \sum_{k=1}^N \log p_{\ell}(x_k | \theta_{\ell}) p(\ell | x_k, \theta^{(i-1)}) = -2 \sum_{k=1}^N \frac{(x_k - \mu_m)}{2\sigma_m^2} p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)}) = 0$$

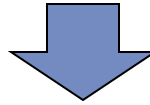


$$\mu_m = \frac{\sum_{k=1}^N x_k p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}{\sum_{k=1}^N p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}$$

moreover, if we let  $c_m = 1/\sigma_m^2$

$$\begin{aligned} \frac{\partial}{\partial c_m} \sum_{\ell=1}^M \sum_{k=1}^N \log p_{\ell}(x_k | \theta_{\ell}) p(\ell | x_k, \theta^{(i-1)}) &= \frac{\partial}{\partial c_m} \sum_{\ell=1}^M \sum_{k=1}^N \left[ -\frac{1}{2} \ln(2\pi\sigma_{\ell}^2) - \frac{(x_k - \mu_{\ell})^2}{2\sigma_{\ell}^2} \right] p(\ell | x_k, \mu_{\ell}^{(i-1)}, \sigma_{\ell}^{(i-1)}) \\ &= \sum_{k=1}^N \left[ \frac{1}{2c_m} - \frac{1}{2} (x_k - \mu_m)^2 \right] p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)}) \\ &= \sum_{k=1}^N \left[ \frac{\sigma_m^2}{2} - \frac{1}{2} (x_k - \mu_m)^2 \right] p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)}) = 0 \end{aligned}$$

$$\frac{\partial}{\partial c_m} \sum_{\ell=1}^M \sum_{k=1}^N \log p_{\ell}(x_k | \theta_{\ell}) p(\ell | x_k, \theta^{(i-1)}) = \sum_{k=1}^N \left[ \frac{\sigma_m^2}{2} - \frac{1}{2} (x_k - \mu_m)^2 \right] p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)}) = 0$$



$$\sigma_m^2 = \frac{\sum_{k=1}^N (x_k - \mu_m)^2 p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}{\sum_{k=1}^N p(m | x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}$$

Finally we find the following set of recursive formulas, that combine the E and M steps:

$$p_m(x|\mu_m, \sigma_m) = \frac{1}{\sqrt{2\pi\sigma_m^2}} \exp\left(-\frac{(x - \mu_m)^2}{2\sigma_m^2}\right)$$

$$p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)}) = \frac{\alpha_m^{(i-1)} p_m(x_k|\mu_m^{(i-1)}, \sigma_m^{(i-1)})}{\sum_{k=1}^M \alpha_m^{(i-1)} p_m(x_k|\mu_m^{(i-1)}, \sigma_m^{(i-1)})}$$

$$\alpha_m^{(i)} = \frac{1}{N} \sum_{k=1}^N p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})$$

$$\mu_m^{(i)} = \frac{\sum_{k=1}^N x_k p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}{\sum_{k=1}^N p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}$$

$$(\sigma_m^{(i)})^2 = \frac{\sum_{k=1}^N (x_k - \mu_m^{(i)})^2 p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}{\sum_{k=1}^N p(m|x_k, \mu_m^{(i-1)}, \sigma_m^{(i-1)})}$$

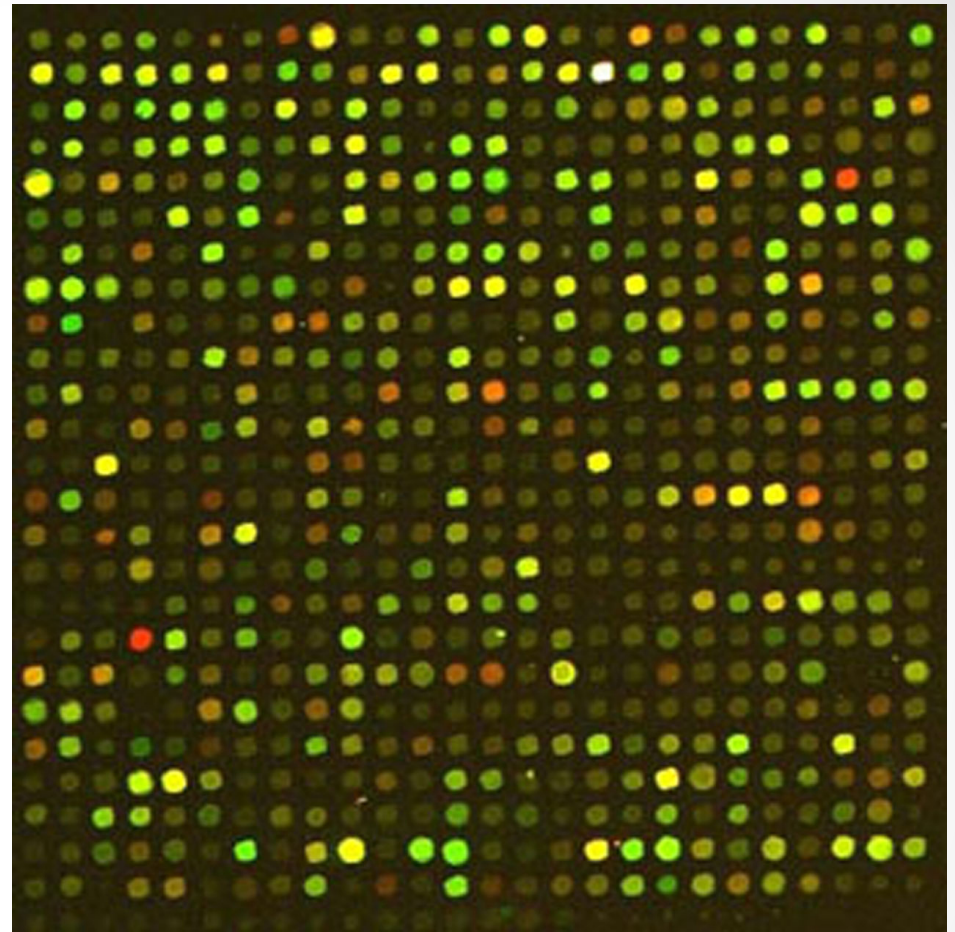
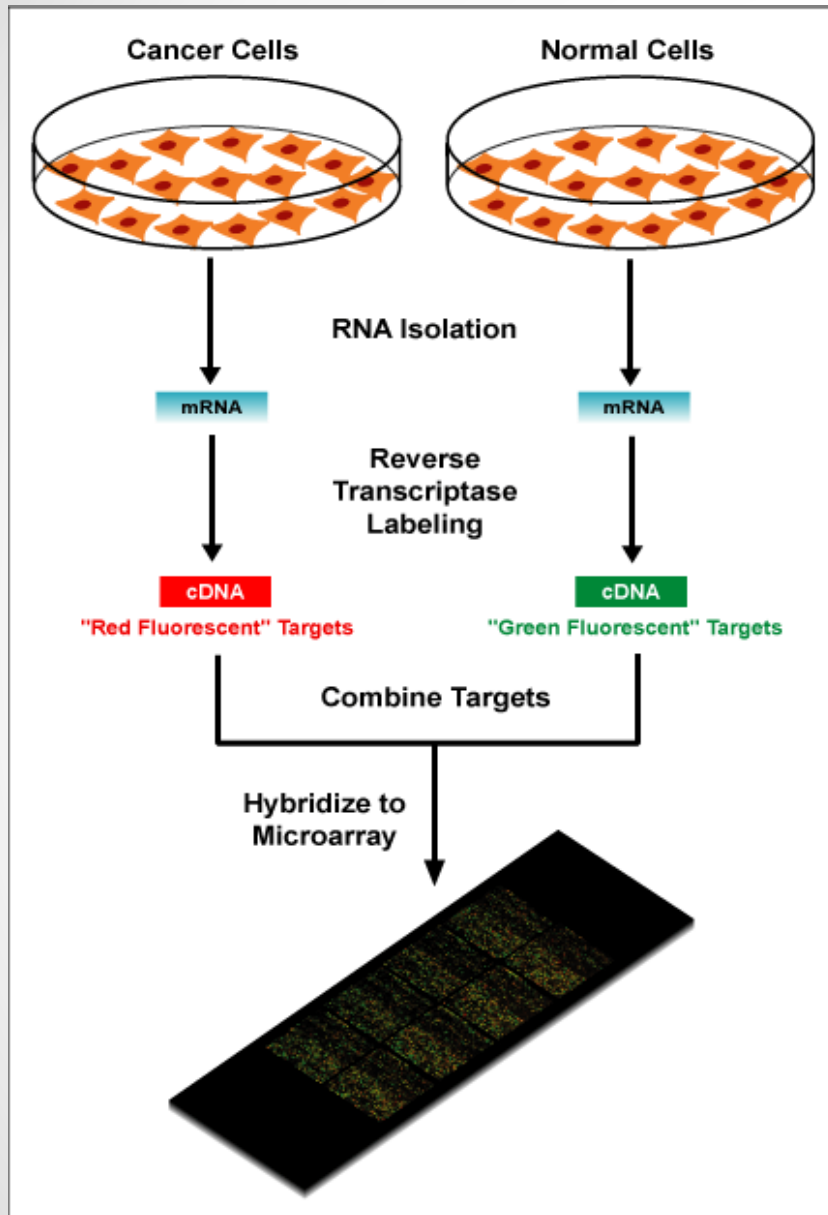
We remark that the probabilities

$$p(y_n | x_n, \boldsymbol{\theta}) = \frac{\alpha_{y_n} p_{y_n}(x_n | \boldsymbol{\theta}_{y_n})}{\sum_{k=1}^M \alpha_k p_k(x_n | \boldsymbol{\theta}_k)}$$

are an estimate of the frequencies of the  $y_n$  using the observed data  $x_n$ , and this is equivalent to a classification (selection of one of the component distributions).



# Example: classification of response of DNA microarrays.

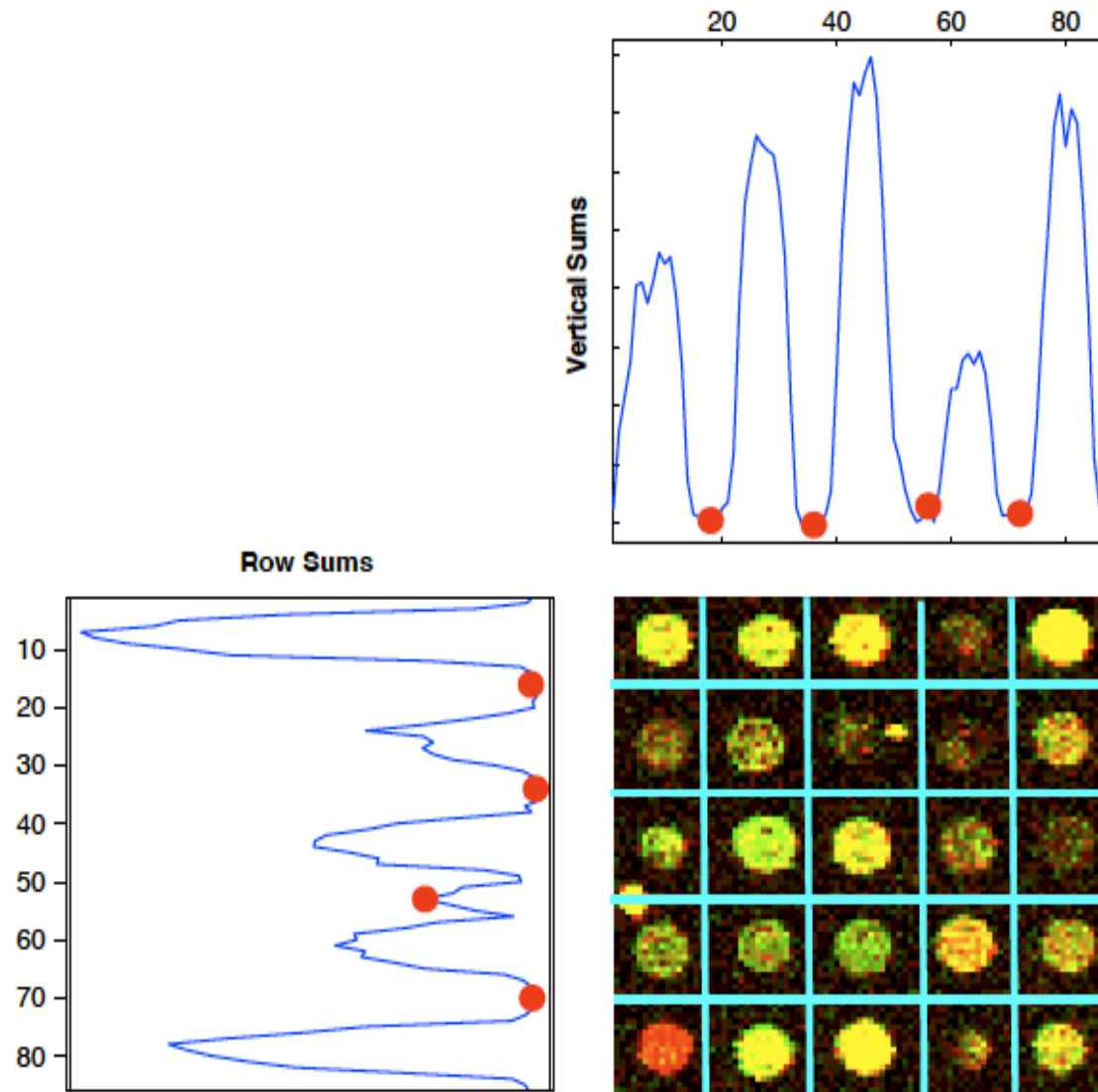


Microarray image from:

[http://www.wormbook.org/chapters/  
www\\_germlinegenomics/germlinegenomics.html](http://www.wormbook.org/chapters/www_germlinegenomics/germlinegenomics.html)

Further informations on DNA microarrays:

[http://www.ncbi.nlm.nih.gov/About/primer/  
microarrays.html](http://www.ncbi.nlm.nih.gov/About/primer/microarrays.html)



From Blekas et al., “Mixture Model Analysis of DNA Microarray Images”, IEEE Trans. on Medical Imaging **24** (2005) 901

## 2. Image processing techniques (MLM, MEM)



The Crab Nebula in Taurus (VLT KUEYEN + FORS2)

ESO PR Photo 40f/99 ( 17 November 1999 )

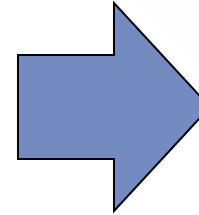
© European Southern Observatory





$f_{11}$	$f_{12}$	$f_{13}$	$f_{14}$	...
$f_{21}$	$f_{22}$	$f_{23}$	$f_{24}$	...
$f_{31}$	$f_{32}$	$f_{33}$	$f_{34}$	...
$f_{41}$	$f_{42}$	$f_{43}$	$f_{44}$	...

pixel map



true  
pixel vector  
**f**

*posterior pixel  
distribution*

*likelihood*

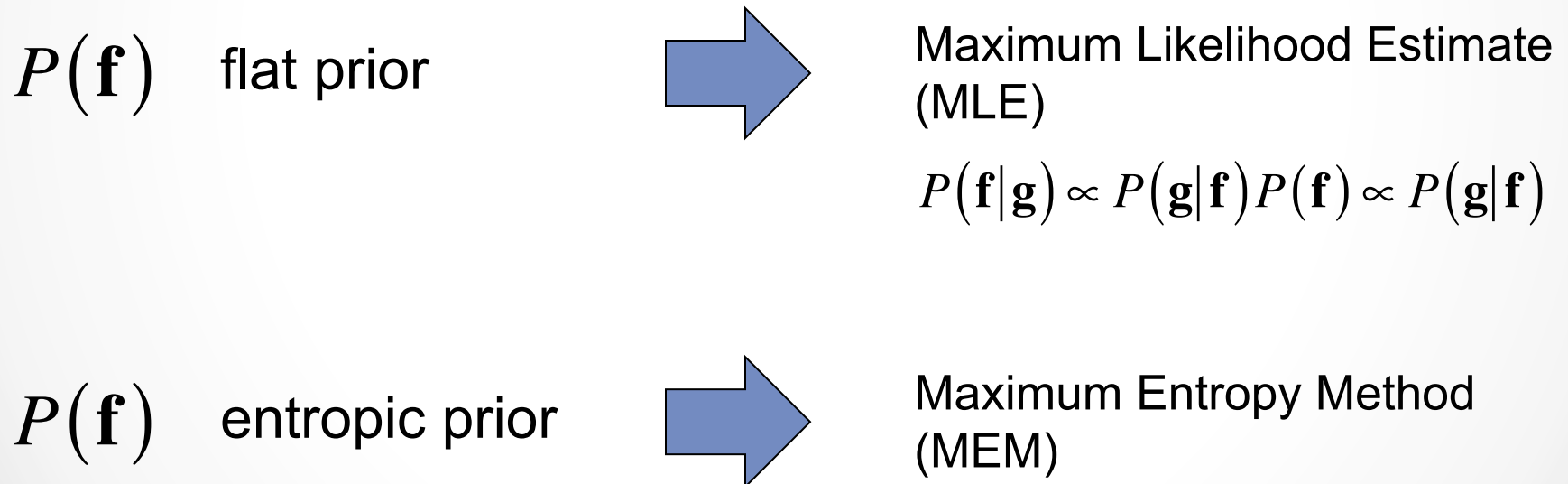
*a priori pixel  
distribution*

$$P(\mathbf{f}|\mathbf{g}) = \frac{P(\mathbf{g}|\mathbf{f})}{P(\mathbf{g})} P(\mathbf{f}) \propto P(\mathbf{g}|\mathbf{f}) P(\mathbf{f})$$

Bayesian estimate of  
true pixel vector from  
observed pixel vector

We estimate the true pixel distribution taking the pixel vector that maximizes the posterior distribution (MAP estimate: Maximum A Posteriori estimate).

This depends on the prior distribution



Notice that

$$\log P(\mathbf{f}|\mathbf{g}) \approx \log P(\mathbf{g}|\mathbf{f}) - [-\log P(\mathbf{f})]$$

therefore we obtain the estimate  $\hat{\mathbf{f}}$  by maximizing the likelihood with the *penalty function*

$$[-\log P(\mathbf{f})]$$

Experiments have been tried with many different penalties, many of them barely justified on probabilistic grounds (or not at all!)

Let  $\mathbf{f}$  be the vector of “true values” (uncorrupted intensities of an image, a spectrum, etc. ...), and translate these values into counts

$$n_i = \lfloor \alpha f_i \rfloor$$

( $i = 1, \dots, M$ ). The least informative prior is that for a structureless image is uniform, and the probability of one count at the  $i$ -th position is just  $1/M$ .

Likewise, the probability of a given vector of values where the total number of counts is  $N$ , is given by the multinomial probability

$$P(\mathbf{n}) = \frac{N!}{n_1! n_2! \dots n_M!} \left( \frac{1}{M} \right)^N ; \quad \sum_k n_k = N$$

## Using Stirling's approximation

$$n! \approx n^n e^{-n} \quad \ln n! \approx n \ln n - n$$

we find, with the definition  $p_i = f_i / \sum_{k=1}^M f_k$

$$\ln P(\mathbf{n}) \approx (N \ln N - N) - \sum_{i=1}^M (n_i \ln n_i - n_i)$$

$$= N \ln N - \sum_{i=1}^M n_i \ln n_i$$

$$\approx -\alpha \sum_{i=1}^M f_i \ln f_i + \text{cost}.$$

*entropic prior*

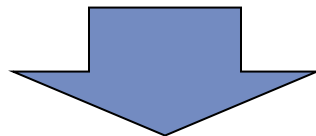


$$P(\mathbf{n}) \propto \exp \left[ -\alpha \sum_{i=1}^M f_i \ln f_i \right] \propto \exp \left[ -\alpha \sum_{i=1}^M p_i \ln p_i \right] = \exp [\alpha S(\mathbf{f})]$$



Using the entropic prior and Bayes' theorem we find

$$P(\mathbf{f}) \propto \exp[\alpha S(\mathbf{f})]$$

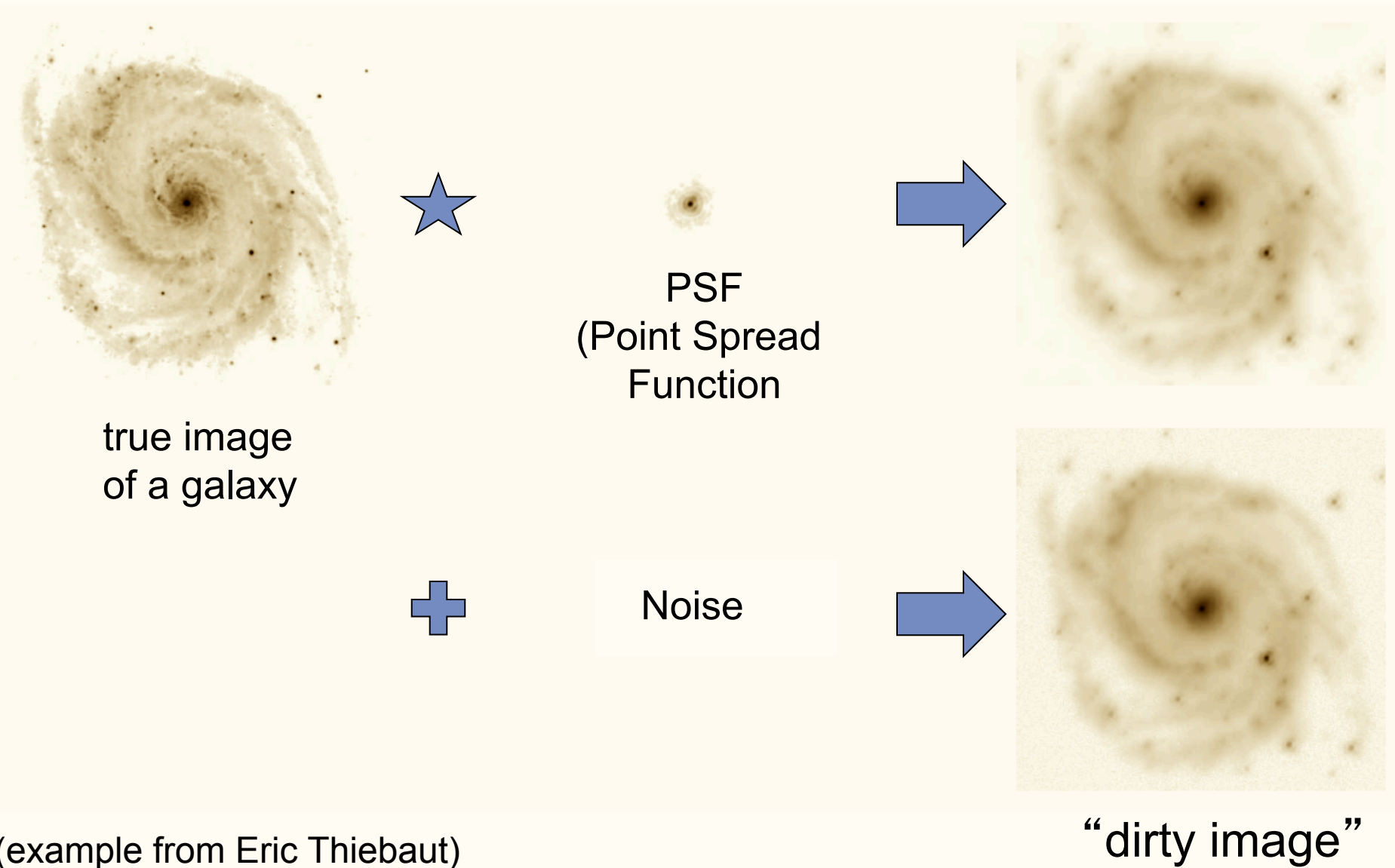


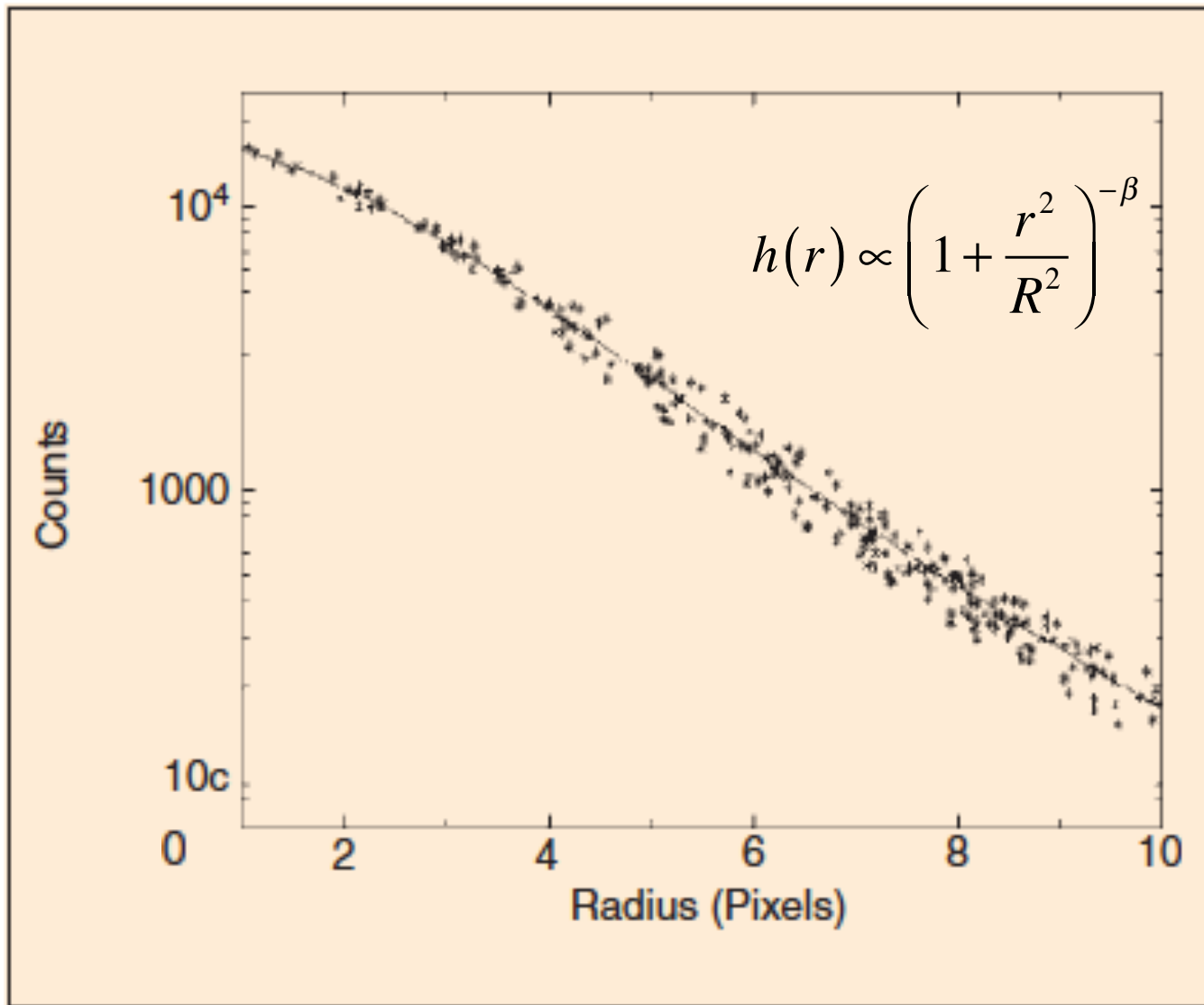
$$P(\mathbf{f}|\mathbf{g}) \propto P(\mathbf{g}|\mathbf{f})P(\mathbf{f}) \propto P(\mathbf{g}|\mathbf{f})\exp[\alpha S(\mathbf{f})]$$

$$\log P(\mathbf{f}|\mathbf{g}) \approx \log P(\mathbf{g}|\mathbf{f}) + \alpha S(\mathbf{f})$$

therefore we find the combination of pixels (i.e., the  $\mathbf{f}$  vector) that maximizes the posterior distribution by maximizing a linear combination of likelihood and Shannon's entropy.

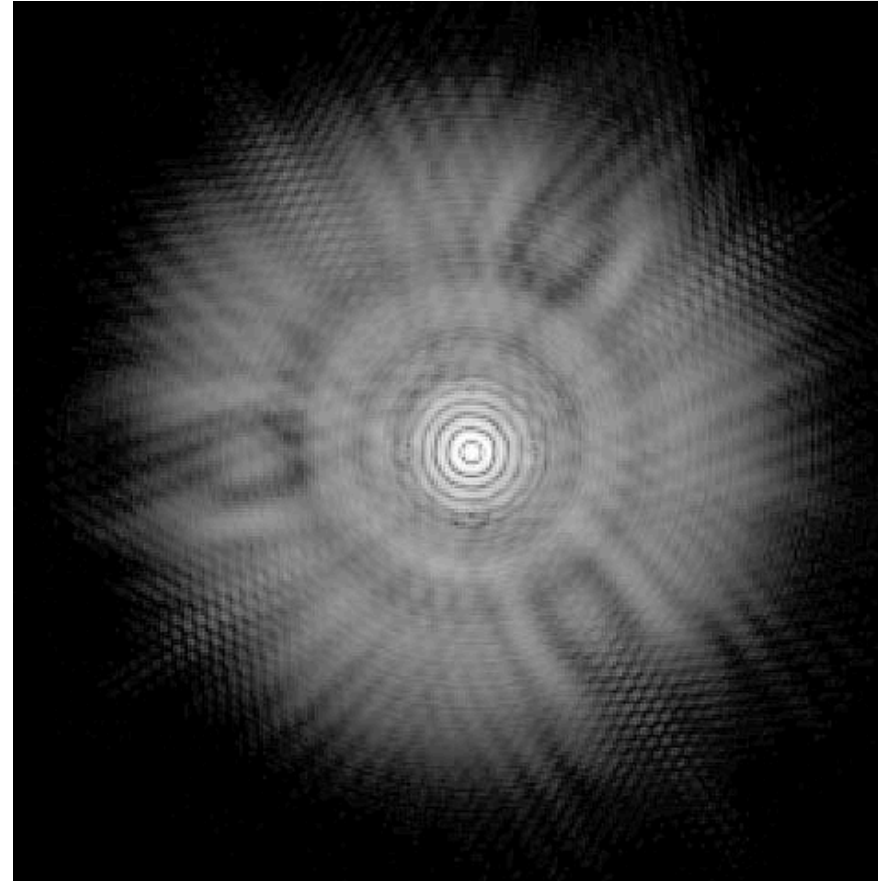
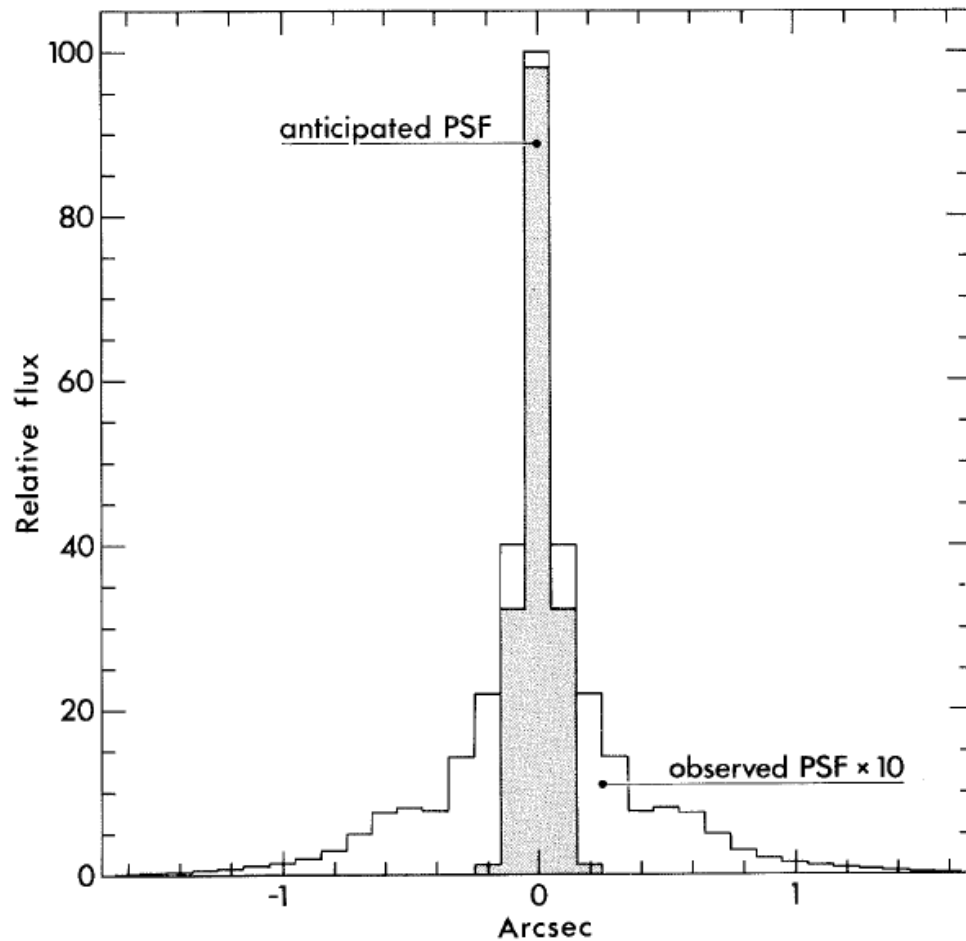
# Image likelihood: 1. the observation model





## PSF from atmospheric turbulence

# The Hubble PSF before the first servicing mission



In general the effect of the PSF is modeled by a linear operator

$$\mathbf{f} \rightarrow \mathbf{H}\mathbf{f}$$

action of optical system  
on true image  
is modeled by matrix **H**

“true” pixel  
vector

## Image likelihood: 2. the noise model (degradation model)

Gaussian noise model

$$P(\mathbf{g}|\mathbf{f}) \propto \exp\left[-\frac{(\mathbf{g} - \mathbf{Hf})^2}{\sigma^2}\right]$$

Poisson noise model

$$P(\mathbf{g}|\mathbf{f}) \propto \prod_n \frac{(\mathbf{Hf})_n^{g_n}}{g_n!} \exp[-(\mathbf{Hf})_n]$$

(Poisson noise mostly from detection process, Gaussian noise mostly from electronics or from approximation of Poisson noise)

sometimes we can use the Gaussian approximation of Poisson noise

$$\begin{aligned} P(\mathbf{g}|\mathbf{f}) &\propto \prod_n \frac{(\mathbf{H}\mathbf{f})_n^{g_n}}{g_n!} \exp[-(\mathbf{H}\mathbf{f})_n] \\ &\approx \prod_n \exp\left[-\frac{(g_n - (\mathbf{H}\mathbf{f})_n)^2}{2(\mathbf{H}\mathbf{f})_n}\right] \\ &= \exp\left[-\sum_n \frac{(g_n - (\mathbf{H}\mathbf{f})_n)^2}{2(\mathbf{H}\mathbf{f})_n}\right] \end{aligned}$$

Gaussian noise only:

maximize linear combination of entropy and chi-square

$$\begin{aligned}\log P(\mathbf{f}|\mathbf{g}) &\approx \alpha S(\mathbf{f}) - \frac{(\mathbf{g} - \mathbf{H}\mathbf{f})^2}{\sigma^2} \\ &= \alpha S(\mathbf{f}) - \sum_n \frac{(g_n - (\mathbf{H}\mathbf{f})_n)^2}{\sigma^2} \\ &= \alpha S(\mathbf{f}) - \chi^2(\mathbf{f})\end{aligned}$$

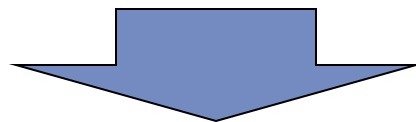


# Combined noise model

detector noise: Poisson noise

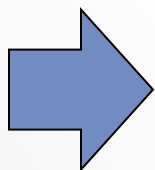
electronic noise: Gaussian noise

$$P(\mathbf{g}|\mathbf{f}) = \prod_n \sum_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(g_n - k)^2}{\sigma^2}\right] \frac{(\mathbf{H}\mathbf{f})_n^k}{k!} \exp[-(\mathbf{H}\mathbf{f})_n]$$



maximize

$$\log P(\mathbf{f}|\mathbf{g}) = \alpha S(\mathbf{f}) + \sum_n \log \left\{ \sum_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(g_n - k)^2}{\sigma^2}\right] \frac{(\mathbf{H}\mathbf{f})_n^k}{k!} \exp[-(\mathbf{H}\mathbf{f})_n] \right\}$$



numerical maximization procedure

# Applications of Max.Ent. to image processing

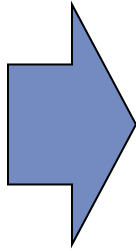
(J. Skilling , Nature **309** (1984) 748)



Car movement introduces linear correlations among pixels. The model of linear corrections does not allow direct inversion to find the corrected image because the number of variables is larger than the number of equations. The MaxEnt methods regularizes the problem and finds a reasonable solution.

# Reconstruction of missing data (from <http://www.maxent.co.uk> )

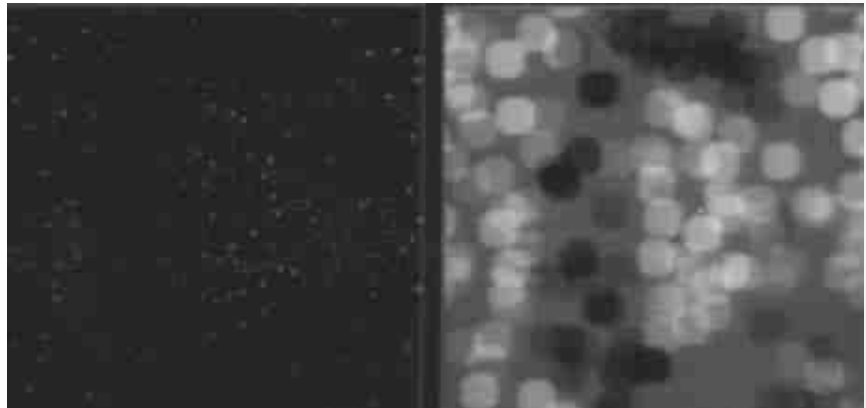
50%

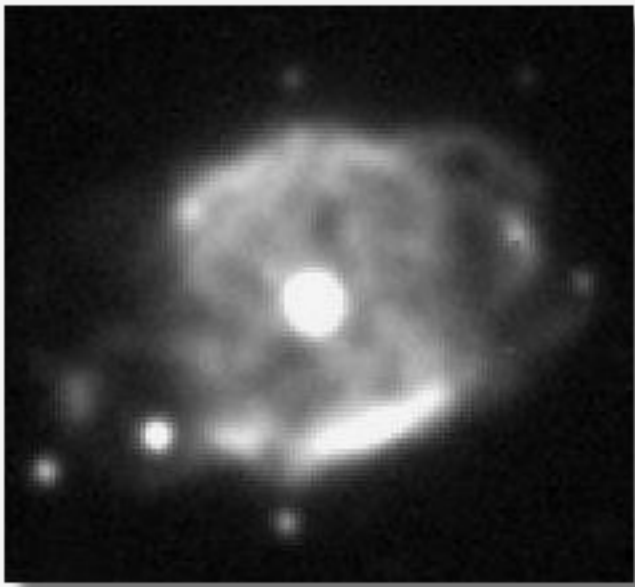


95%

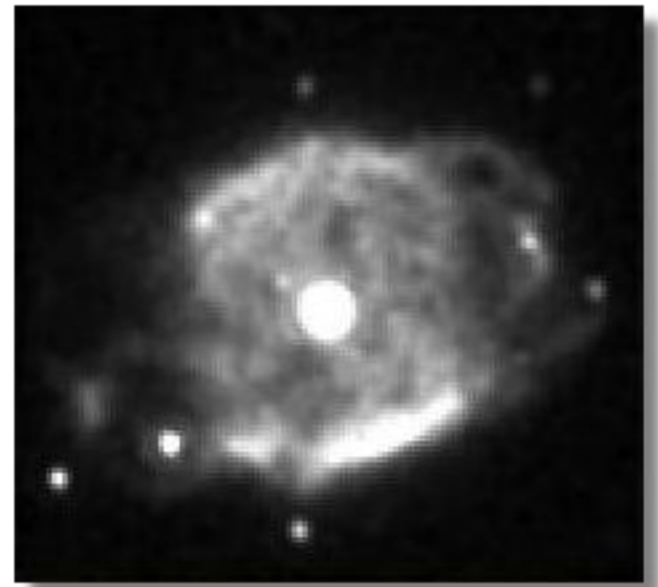


99%





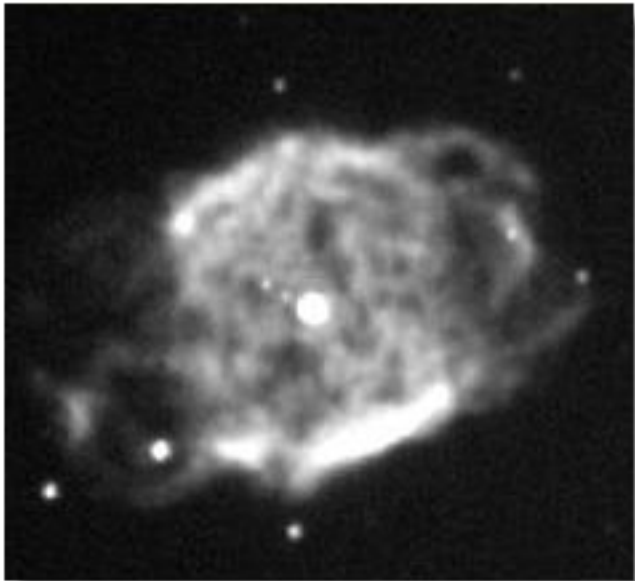
NGC 40



low resolution (MEM enhanced)

low resolution

high resolution



### John Skilling: Biographical information

John is Scientific Director of MEDC. He did his Ph.D. (on cosmic rays) in the Department of Physics at Cambridge University, and went on to become a Lecturer in the Department of Applied Mathematics and Theoretical Physics, and a Fellow of St Johns College.

In the late 1970s, another radio astronomer, [Steve Gull](#), introduced him to the power of the Maximum Entropy Method. John wrote what was to become the first MemSys kernel system, and helped lay the Bayesian foundations for MEM. In 1981 he and Steve founded MEDC to exploit opportunities to apply MEM in other fields.

John resigned his Lectureship in 1990 in order to go fulltime with MSL and MEDC. Thanks to the wonders of modern technology John is able to telecommute from his new home in the West of Ireland, and he makes regular visits to clients both in the UK and further afield.



[Home](#) | [Applications](#) | [Products](#) | [Prices](#) | [Documents](#) | [About MEDC](#) | [Contact Us](#) | [Full search](#)

[Home](#)

[About MEDC](#)

[Applications](#)

[Examples](#)

[Products](#)

[Prices](#)

[Documents](#)

[Contact us](#)

[Search MEDC](#)

Quick Search:

Search

©MEDC Ltd. Last revised Wed Sep 19 22:19:39 2007

<http://www.maxent.co.uk/>

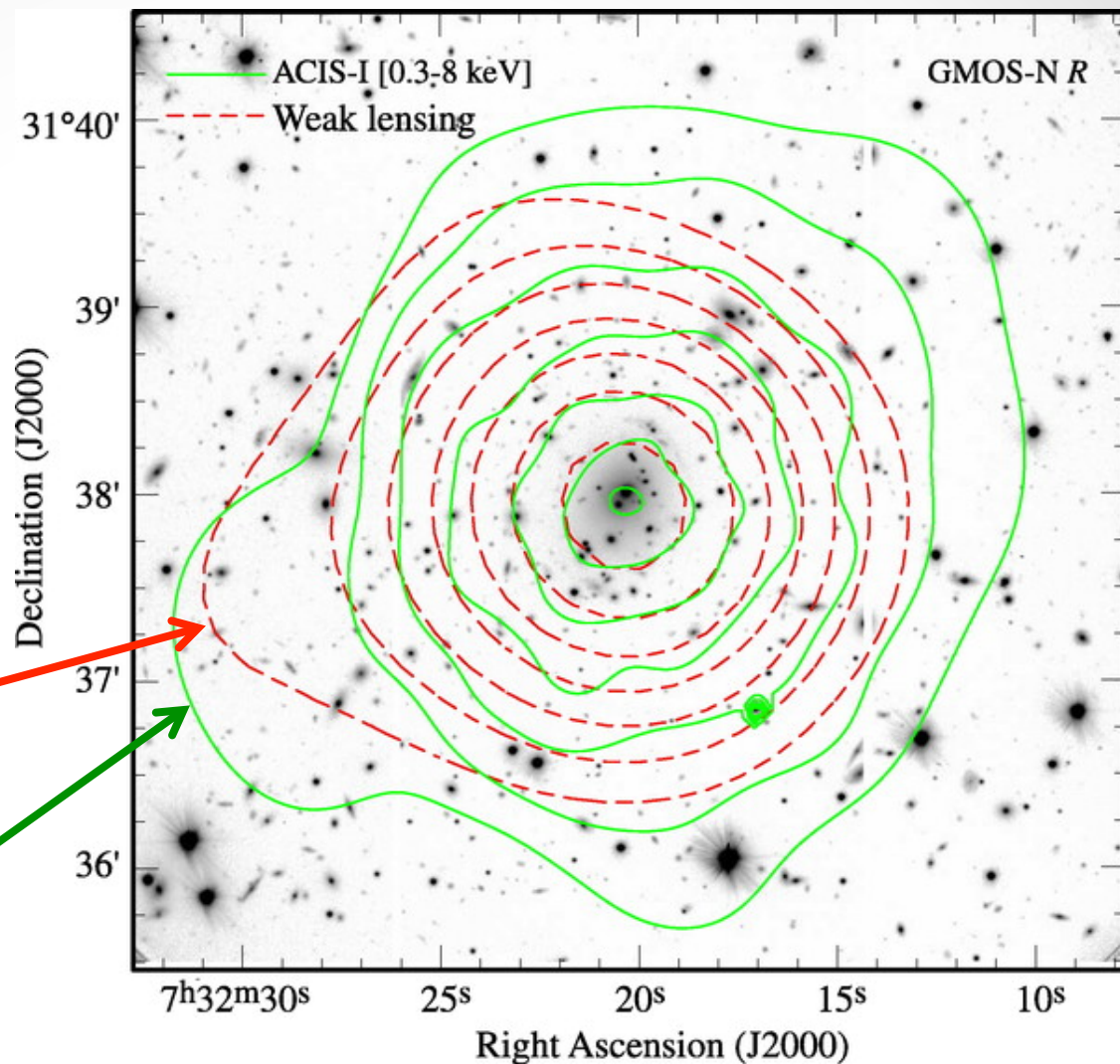


Example of LensEnt usage (Bridle et al, 1998):

reconstruction of mass density from lensing data, using Max Ent

reconstructed mass density

X-ray emission data



GMOS image of the central region of Abell 586 with **logarithmically spaced X-ray isophotes (solid lines) and weak-lensing reconstructed mass density (dashed lines) superposed**. The X-ray point source near the southwest corner is the Seyfert 1 galaxy C171\_3650.  
(from Cypriano et al., ApJ, **630** (2005) 38-49)

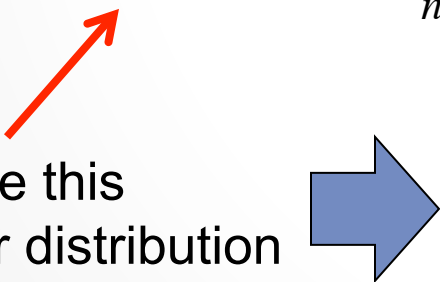
Many related methods: e.g. the Richardson-Lucy (RL) algorithm

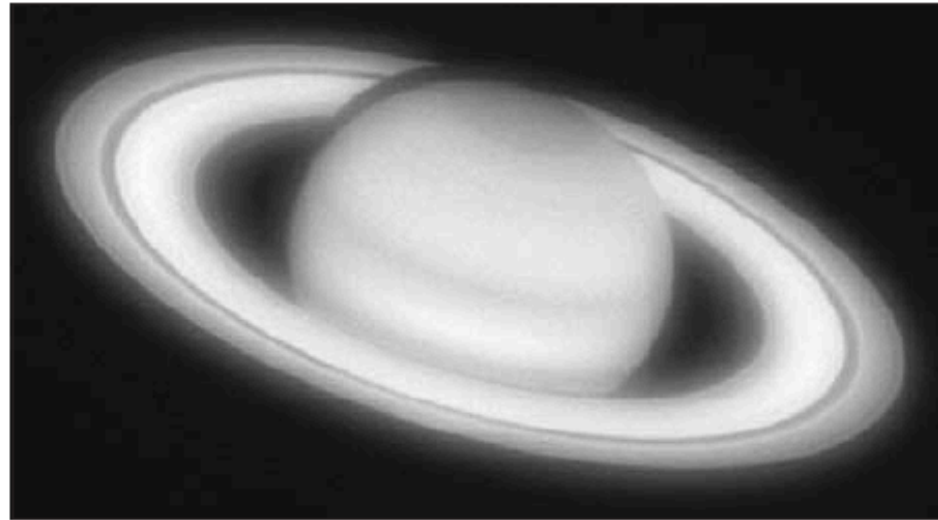
noise model: Poisson noise

prior: flat prior

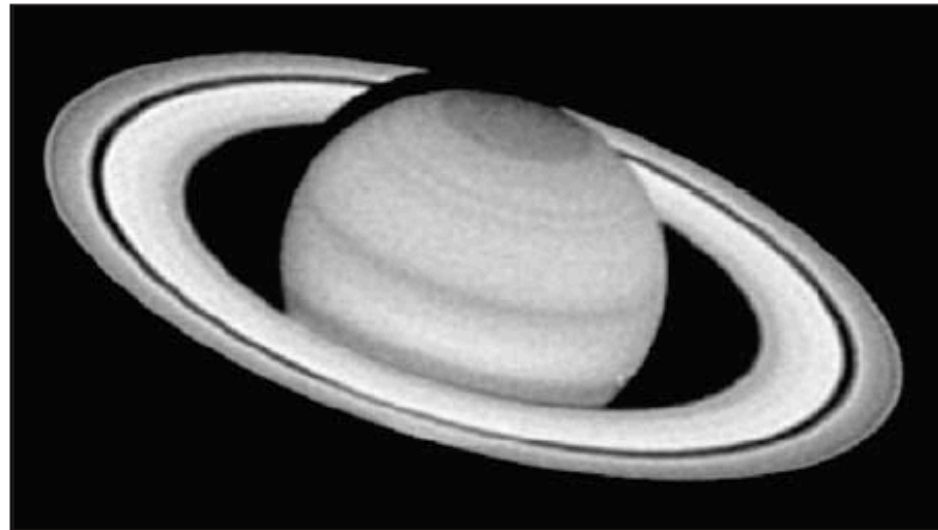
$$P(\mathbf{f}|\mathbf{g}) \propto \prod_n \frac{(\mathbf{Hf})_n^{g_n}}{g_n!} \exp[-(\mathbf{Hf})_n] P(\mathbf{f})$$

$$\log P(\mathbf{f}|\mathbf{g}) \approx \sum_n \left[ -(\mathbf{Hf})_n + g_n \log(\mathbf{Hf})_n \right] + \text{const.}$$

maximize this  
posterior distribution   $\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \sum_n \left[ -(\mathbf{Hf})_n + g_n \log(\mathbf{Hf})_n \right]$



▲ 8. *Raw image of planet Saturn obtained with the WF/PC camera of the HST.*



▲ 9. *Reconstruction of the image of Saturn using the R-L algorithm.*



## NGC 604 in Spiral Galaxy M33



# References:

## EM algorithm:

- A. P. Dempster, N. M. Laird, and D. B. Rubin: “Maximum likelihood from incomplete data via the EM algorithm”, Journal of the Royal Statistical Society series B, **39** (1977) 1
- J. Bilmes: “A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models”, ICSI preprint TR-97-021 (1998)
- B. Flury and A. Zoppé, “Exercises in EM” American Statistician, **54** (2000) 207.

## MaxEnt and image processing

- J. Skilling et al., Mon. Not. R. astr. Soc. **187** (1979) 145
- J. Skilling , Nature **309** (1984) 748
- R. Narayan and R. Nityananda, Ann. Rev. Astron. Astrophys. **24** (1986) 127
- R. Molina et al., IEEE Signal Proc. Magazine (marzo 2001) 13
- J. Skilling, A. W. Strong and K. Bennett, Mon. Not. R. astr. Soc. **187** (1979) 145
- J. Skilling and R. K. Bryan, Mon. Not. R. astr. Soc. **211** (1984) 11
- S. L. Bridle et al, Mon. Not. R. astr. Soc. **299** (1998) 895