

Some Model Selection Problems in Image Analysis

Mário A. T. Figueiredo

“Instituto Superior Técnico”

Lisboa

PORTUGAL

People involved in work here presented:

José M. N. Leitão, Instituto Superior Técnico, PORTUGAL

Anil K. Jain, Michigan State University, USA

Robert D. Nowak, Rice University, USA

2 Outline

- Image analysis as Bayesian inference: a brief review
- Model selection problems
- Approaches to model selection:

Bayesian model selection (Bayes factors)

Minimum description length (MDL)

- Example: discontinuity-preserving restoration
- Example: contour estimation
- Example: segmentation of Poisson imagery
- Concluding remarks

3 Image analysis as statistical inference

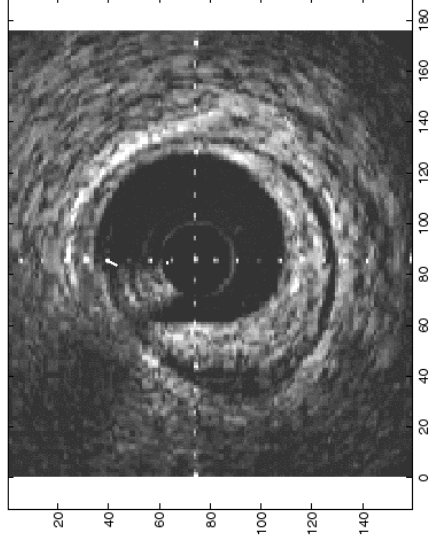
Most image analysis problems can/should be formulated as

Given observed data \mathbf{g} , infer \mathbf{f}

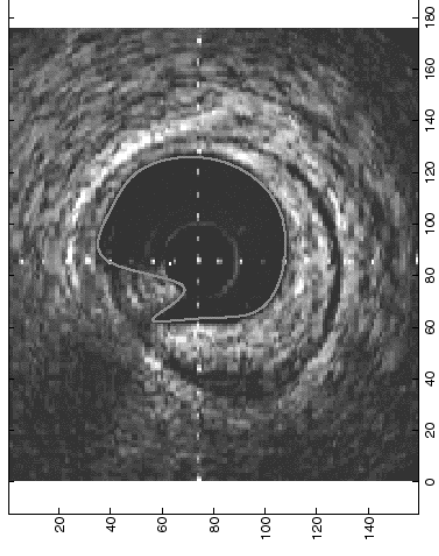
Seems a trivial statement.

My message: “start by writing down what \mathbf{f} and \mathbf{g} are”

- Examples



inference



\mathbf{g} , an observed image

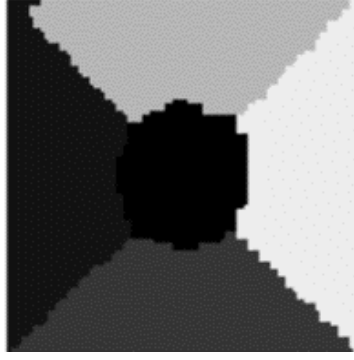
\mathbf{f} , a contour, e.g., represented by a sequence of points

4 Image analysis as statistical inference: more examples

g, observed image



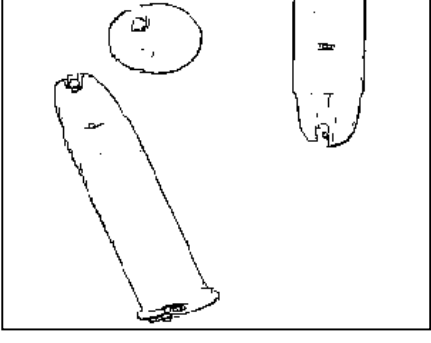
inference
↓



f, segmentation



inference
↑



g, observed image

f, edge map



inference
↑



g, noisy image

f, restored image

5 Image analysis as statistical inference

- In all the cases here considered: \mathbf{g} = “one observed image”
- Other possibilities:
 - extracted features (e.g., local filters),
 - tomographic projections,
 - sequences of images,
 -
- \mathbf{f} can be image-like (e.g., in image restoration or segmentation)
- \mathbf{f} can be a parametric representation (e.g., a spline contour)

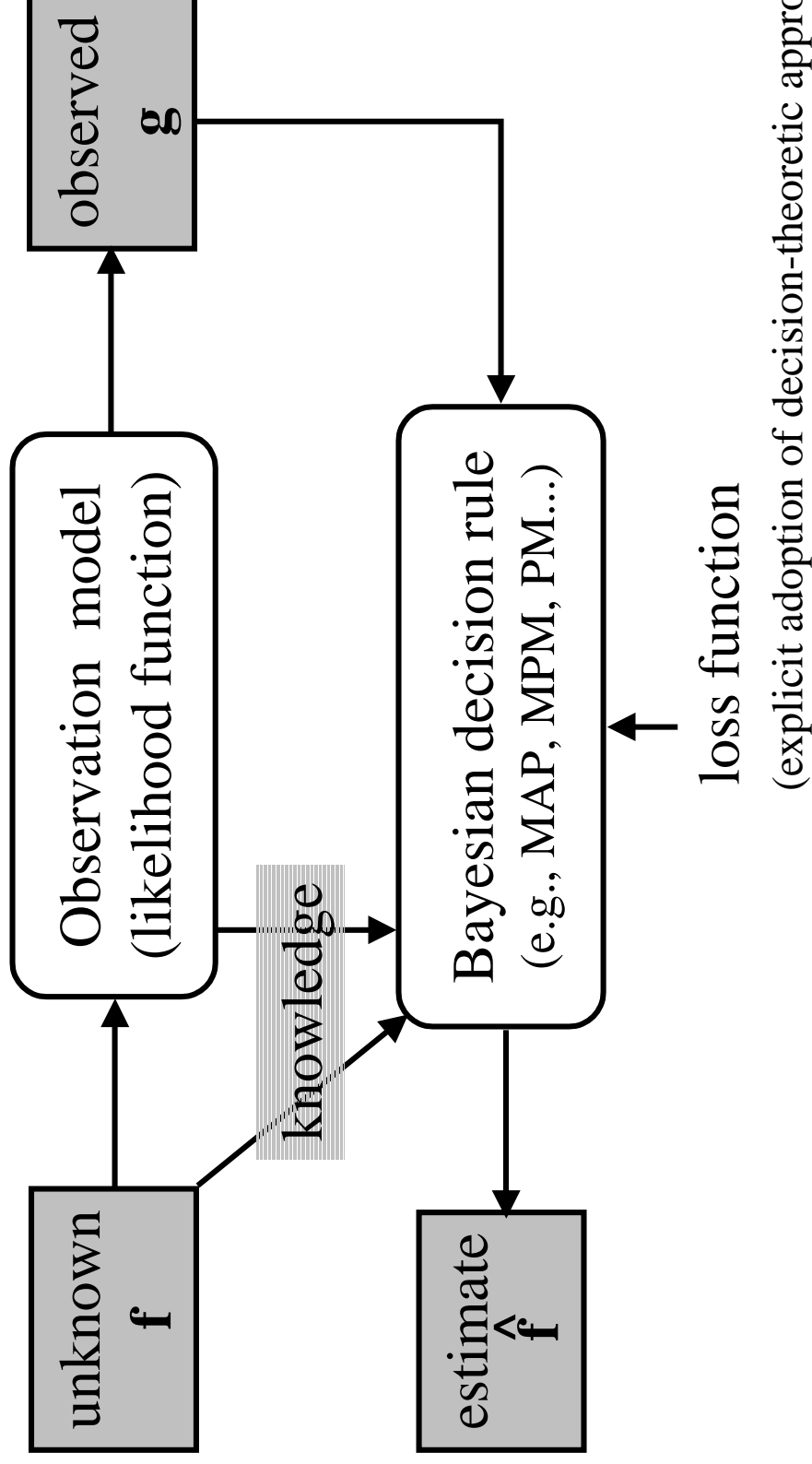
Many uncertainties
(noise, uncertainty about \mathbf{f}, \dots)



Natural frameworks:

- statistical inference
- statistical decision making

6 Bayesian image analysis



The Bayesian approach is explicitly model-based

7 Bayesian image analysis: basic elements

- Observation model / likelihood function:

$$p(\mathbf{g} | \mathbf{f}, \phi)$$

\mathbf{f} is the unknown

\mathbf{g} is the observed data

ϕ are parameters (assumed known, by now)

- Prior knowledge:

$$p(\mathbf{f} | \psi)$$

\mathbf{f} is the unknown

ψ are parameters (assumed known, by now)

- *A posteriori* knowledge (via Bayes theorem):

$$p(\mathbf{f} | \mathbf{g}, \phi, \psi) = \frac{p(\mathbf{g} | \mathbf{f}, \phi) p(\mathbf{f} | \psi)}{p(\mathbf{g} | \phi, \psi)}$$

8 Bayesian image analysis: loss function and Bayes risk

Given $p(\mathbf{f} | \mathbf{g}, \phi, \psi)$ (obtained via Bayes law)

...and a loss function $L(\mathbf{f}, \hat{\mathbf{f}})$

A posteriori expected loss:

$$E[L(\mathbf{f}, \hat{\mathbf{f}}) | \mathbf{g}, \phi, \psi] = \int L(\mathbf{f}, \hat{\mathbf{f}}) p(\mathbf{f} | \mathbf{g}, \phi, \psi) d\mathbf{f}$$

...expected value of the loss function, given the observed data \mathbf{g} .

Optimal Bayes rule

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}'} E[L(\mathbf{f}, \mathbf{f}') | \mathbf{g}, \phi, \psi]$$

...minimizes the *a posteriori* expected loss:

9 Bayesian image analysis: common decision rules

So, what to do with the posterior $p(\mathbf{f} | \mathbf{g}, \phi, \psi)$?

It all depends on the loss function. Examples:

$$L(\mathbf{f}, \hat{\mathbf{f}}) = \begin{cases} 1 & \Leftarrow \mathbf{f} \neq \hat{\mathbf{f}} \\ 0 & \Leftarrow \mathbf{f} = \hat{\mathbf{f}} \end{cases} \longrightarrow \hat{\mathbf{f}}_{\text{MAP}} = \arg \max_{\mathbf{f}} p(\mathbf{f} | \mathbf{g}, \phi, \psi)$$

$$L(\mathbf{f}, \hat{\mathbf{f}}) = \|\mathbf{f} - \hat{\mathbf{f}}\|^2 \longrightarrow \hat{\mathbf{f}}_{\text{PM}} = E[\mathbf{f} | \mathbf{g}, \phi, \psi]$$

$$L(\mathbf{f}, \hat{\mathbf{f}}) = \sum_i L_i(\mathbf{f}_i, \hat{\mathbf{f}}_i) \longrightarrow \hat{\mathbf{f}}_i = \arg \min_{\mathbf{f}_i} E[L_i(\mathbf{f}_i - \hat{\mathbf{f}}_i) | \mathbf{g}, \phi, \psi]$$

(additive loss)

Marginal criteria, e.g., *maximizer of posterior marginals* (MPM), Marroquin, Mitter, and Poggio, JASA 1987)

10 Unknown (hyper)parameters

What if ϕ and/or ψ (the hyper-parameters) are unknown ?

- Unknown: \mathbf{f}, ϕ, ψ
- *A posteriori* probability function

$$p(\mathbf{f}, \phi, \psi | \mathbf{g}) = \frac{p(\mathbf{g} | \mathbf{f}, \phi) p(\mathbf{f} | \psi) p(\phi, \psi)}{p(\mathbf{g})}$$

$p(\phi, \psi)$
“hyper-prior”

- There is nothing fundamentally new
- This is called a hierarchical Bayes model

Again, it all depends on the loss function. Examples:

$$L(\mathbf{f}, \phi, \psi), (\hat{\mathbf{f}}, \hat{\phi}, \hat{\psi}) = \begin{cases} 1 & \Leftarrow \mathbf{f} \neq \hat{\mathbf{f}} \\ 0 & \Leftarrow \mathbf{f} = \hat{\mathbf{f}} \end{cases} \longrightarrow \hat{\mathbf{f}}_{\text{MMAP}} = \arg \max_{\mathbf{f}} p(\mathbf{f} | \mathbf{g})$$

(*marginal MAP*)

(does not care about ϕ and ψ)

$$L(\cdot, \cdot) = \begin{cases} 1 & \Leftarrow (\mathbf{f}, \phi, \psi) \neq (\hat{\mathbf{f}}, \hat{\phi}, \hat{\psi}) \\ 0 & \Leftarrow (\mathbf{f}, \phi, \psi) = (\hat{\mathbf{f}}, \hat{\phi}, \hat{\psi}) \end{cases} \longrightarrow (\hat{\mathbf{f}}, \hat{\phi}, \hat{\psi})_{\text{GMAP}} = \arg \max_{\mathbf{f}, \phi, \psi} p(\mathbf{f}, \phi, \psi | \mathbf{g})$$

(*generalized MAP*)

$$L(\cdot, \cdot) = \begin{cases} 1 & \Leftarrow (\phi, \psi) \neq (\hat{\phi}, \hat{\psi}) \\ 0 & \Leftarrow (\phi, \psi) = (\hat{\phi}, \hat{\psi}) \end{cases} \longrightarrow (\hat{\phi}, \hat{\psi})_{\text{PMMAP}} = \arg \max_{\phi, \psi} p(\phi, \psi | \mathbf{g})$$

(*parameter marginal MAP, usually via EM*)

$$\text{If } p(\phi, \psi) \propto \text{const} \quad (\hat{\phi}, \hat{\psi})_{\text{PMMAP}} = (\hat{\phi}, \hat{\psi})_{\text{MML}} = \arg \max_{\phi, \psi} p(\mathbf{g} | \phi, \psi)$$

(maximum marginal likelihood)

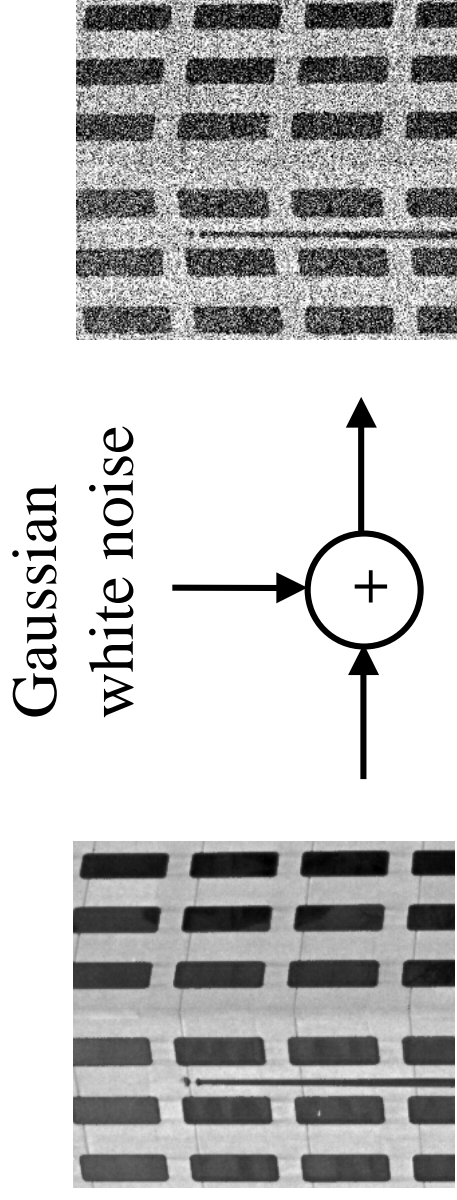
12 Empirical Bayes

1. Estimate hyper-parameters, e.g., $(\hat{\phi}, \hat{\psi})_{\text{MML}} = \arg \max_{\phi, \psi} p(\mathbf{g} | \phi, \psi)$
2. Plug them in as if they were known, e.g.

$$\hat{\mathbf{f}}_{\text{MAP}} = \arg \max_{\mathbf{f}} p(\mathbf{f} | \mathbf{g}, \hat{\phi}, \hat{\psi})$$

This is **not** a truly Bayesian procedure,

...but has nice asymptotic properties.



Likelihood function

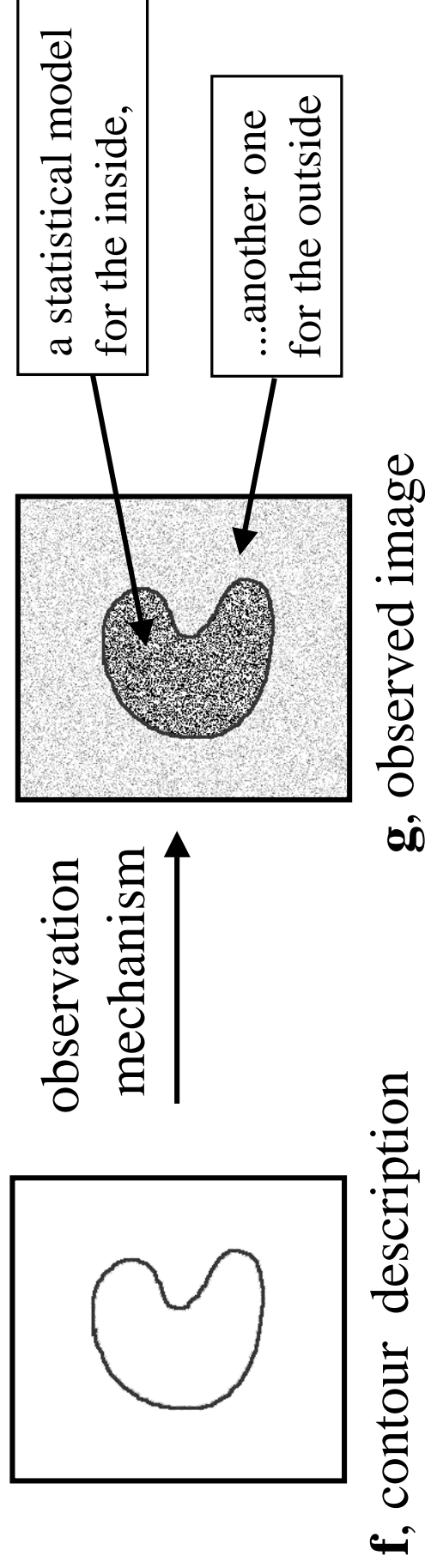
$$p(\mathbf{g}|\mathbf{f}, \sigma^2) \equiv \mathcal{N}(\mathbf{g}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$$

$\sigma^2 \rightarrow$ noise variance

$N \rightarrow$ number of pixels

$\mathbf{I}_N \rightarrow$ $N \times N$ identity matrix

See: M. Figueiredo and J. Leitão, “Unsupervised image restoration and Edge location using Gauss-Markov random fields and the MDL principle”, *IEEE Transactions on Image Processing*, vol. 6, pp. 1089-1102, 1997.



Under inside/outside independence assumption:

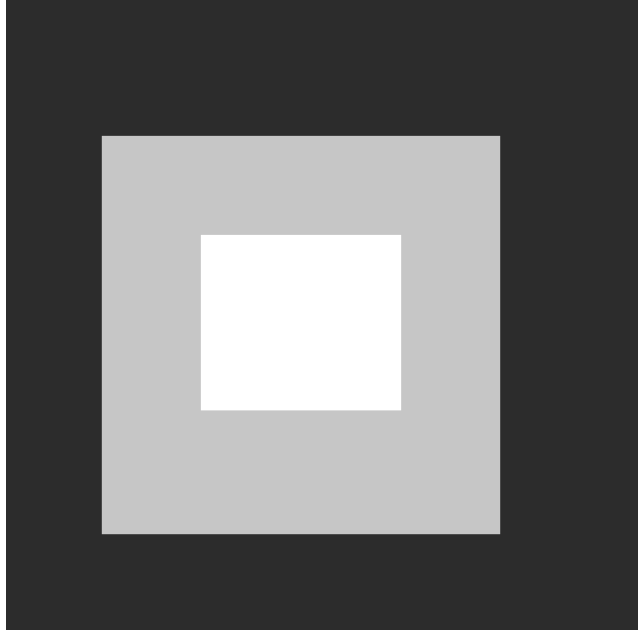
$$p(\mathbf{g} | \mathbf{f}, \phi_{\text{in}}, \phi_{\text{out}}) = p(\{g_i : i \in \text{in}(\mathbf{f})\} | \phi_{\text{in}}) p(\{g_i : i \in \text{out}(\mathbf{f})\} | \phi_{\text{out}})$$

Examples: Gaussians of different means and variances;

Rayleigh of different variances (ultrasound images); Different textures.

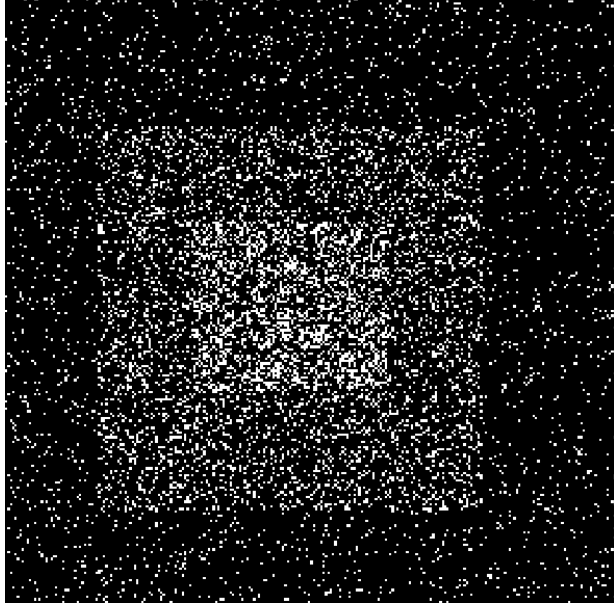
See: M. Figueiredo, J. Leitão, and A. Jain, “Adaptive B-splines and boundary estimation”, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR-97), pp. 724-729, San Juan, Puerto Rico, 1997.

15 Image analysis examples: Poisson image segmentation



f, intensity image

observation
mechanism
→



g, observed counts

Under conditional independence assumption:

$$p(\mathbf{g} | \mathbf{f}) = \prod_i p(g_i | f_i) = \prod_i \frac{(f_i)^{g_i} \exp(-f_i)}{g_i!}$$

(independent
Poisson
observations)

What prior information to use ?

- In image restoration → Images should be smooth, except at edges/discontinuities
- In contour estimation → Contours should be smooth
- In image segmentation → Nearby pixels tend to belong to the same region

Questions: How many edges should we allow ?

How smooth should the contour be ?

Into how many regions should we segment the image ?

These are all model selection questions

17 Model selection

Suppose there are K possible models, $m \in \{m_1, \dots, m_K\}$

...with a priori probabilities $p(m_1), \dots, p(m_K)$

Given model m :

- a likelihood function

$$p(\mathbf{g} | \mathbf{f}, \phi_{(m)}, m)$$

- a prior

$$p(\mathbf{f} | \psi_{(m)}, m)$$

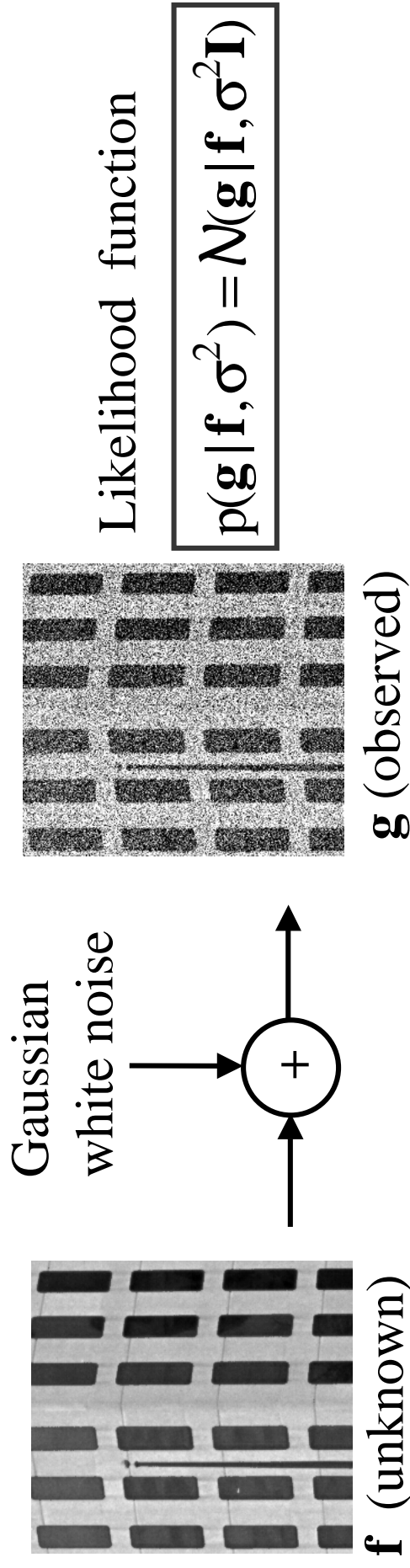
- a hyper-prior

$$p(\phi_{(m)}, \psi_{(m)} | m)$$

$\phi_{(m)}, \psi_{(m)}$ are indexed by m : they may have different structures

There is nothing new: unknowns are now $m, \mathbf{f}, \phi_{(m)}, \psi_{(m)}$
standard Bayes approach can be followed

Model selection: image restoration



Prior (smoothness almost everywhere)

$$p(\mathbf{f} | \boldsymbol{\varphi}_{(m)}, \mathbf{m}) = \frac{1}{Z(\boldsymbol{\varphi}_{(m)}, \mathbf{m})} \exp\left[-\frac{\mu}{2} \sum_{i \sim j} (1 - d_{i,j}) (\mathbf{f}_i - \mathbf{f}_j)^2\right]$$

number of $d_{i,j}=1$ (points to the sum)
 partition function (points to $Z(\boldsymbol{\varphi}_{(m)}, \mathbf{m})$)
 $d_{i,j} \in \{0, 1\}$ (points to $d_{i,j}$)
 sum over all nearest neighbors (points to the sum)
 $\boldsymbol{\varphi}_{(m)} = (\mu, \{\mathbf{1}, \mathbf{j}\} : d_{i,j} = 1)$
 edge locations (points to the set definition)

19 Model selection: image restoration

Compound Gauss-Markov random field (CGMRF)

$$p(\mathbf{f} | \Psi_{(m)}, \mathbf{m}) = \frac{\sqrt{\det \mathbf{A}(\Psi_{(m)})}}{(2\pi)^{N/2}} \exp \left[-\frac{1}{2} \mathbf{f}^T \mathbf{A}(\Psi_{(m)}) \mathbf{f} \right]$$

$\mathbf{A}(\Psi_{(m)}) \rightarrow$ “potential matrix” (inverse of the covariance)

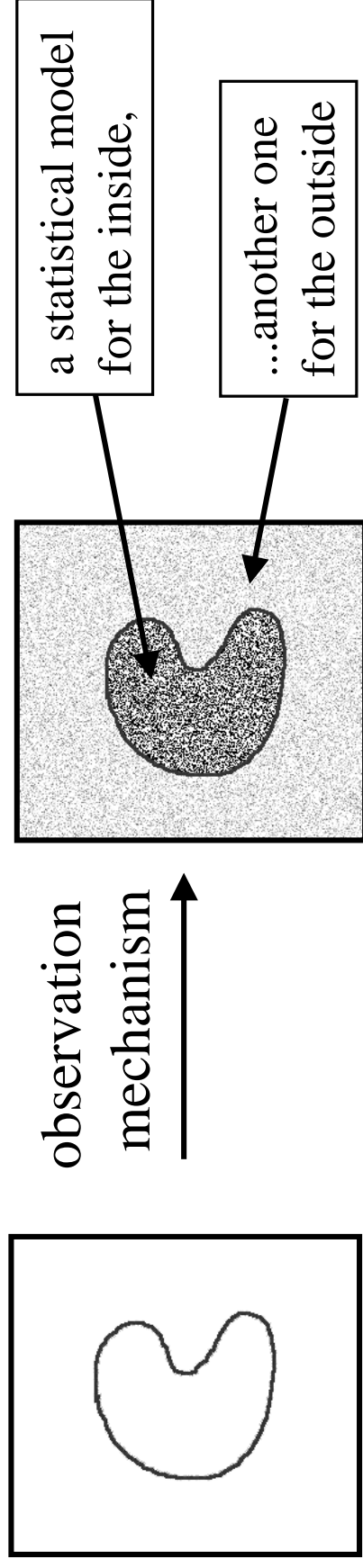
Given $\Psi_{(m)}$, the image estimate is given by a simple Wiener-type filter:

$$\hat{\mathbf{f}} = \left[\sigma^2 \mathbf{A}(\Psi_{(m)}) + \mathbf{I} \right]^{-1} \mathbf{g}$$

But $\mathbf{m} = ?$ How many edges? \rightarrow model selection problem

See: M. Figueiredo and J. Leitão, “Unsupervised image restoration and Edge location using Gauss-Markov random fields and the MDL principle”, *IEEE Transactions on Image Processing*, vol. 6, pp. 1089-1102, 1997.

Model selection: contour estimation



f, contour description

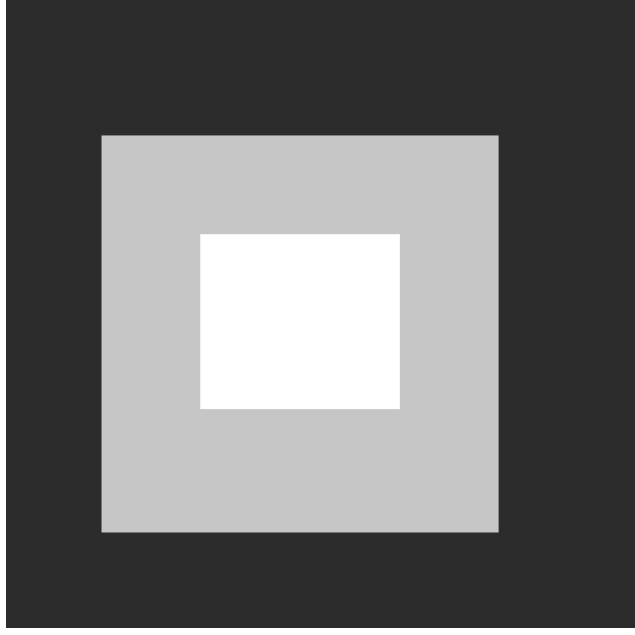
g, observed image

A (discretized) B-spline contour description

$\mathbf{f}_{(m)} = \mathbf{B}_{(m)} \boldsymbol{\theta}_{(m)}$ \rightarrow a 2-D spline curve with m control points

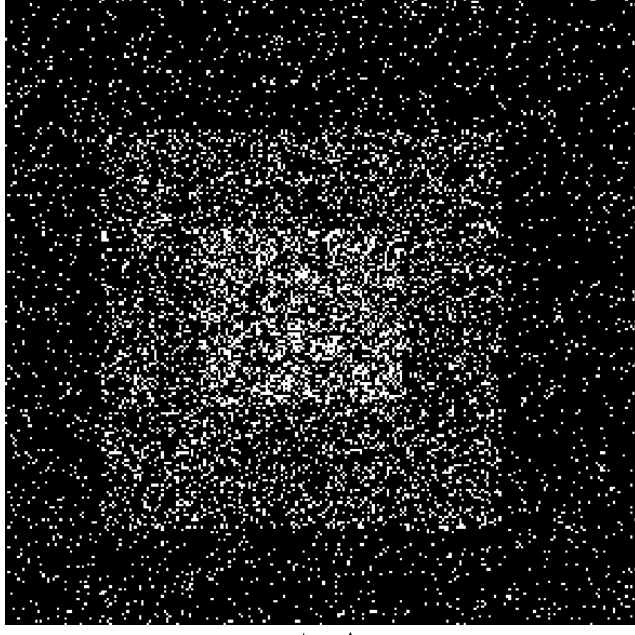
matrix with discretized B-spline basis \rightarrow control points coordinates

$m = ?$ How complex should the contour be? \rightarrow model selection



f , intensity image

observation
mechanism



g , observed counts

f is a piece-wise constant image

How many regions of constant intensity ?



model selection problem

$$p(\mathbf{m}, \mathbf{f}, \phi_{(m)}, \Psi_{(m)} | \mathbf{g}) = \frac{p(\mathbf{g} | \mathbf{f}, \phi_{(m)}, \mathbf{m}) p(\mathbf{f} | \Psi_{(m)}, \mathbf{m}) p(\phi_{(m)}, \Psi_{(m)} | \mathbf{m}) p(\mathbf{m})}{p(\mathbf{g})}$$

Again, what to do with this, depends on the loss function (“true Bayesians” just report it, or sample from it; e.g., RJ-MCMC)

Seen strictly as a model selection problem, the natural choice is

$$L(\cdot, \cdot) = \begin{cases} 1 & \Leftarrow m \neq \hat{m} \\ 0 & \Leftarrow m = \hat{m} \end{cases} \longrightarrow \hat{m} = \arg \max_m p(\mathbf{m} | \mathbf{g})$$

$$\hat{m} = \arg \max_m \{ p(\mathbf{m} | \mathbf{g}) \} = \arg \max_m \{ \underbrace{p(\mathbf{g} | \mathbf{m}) p(\mathbf{m})}_{\text{evidence for model } m} \}$$

sometimes called the
“evidence for model m ”

$$\hat{m} = \arg \max_m \{ p(m | \mathbf{g}) \} = \arg \max_m \{ p(\mathbf{g} | m) p(m) \}$$

When comparing two models

$$\frac{p(m_1 | \mathbf{g})}{p(m_2 | \mathbf{g})} = \frac{p(\mathbf{g} | m_1) p(m_1)}{p(\mathbf{g} | m_2) p(m_2)} \underbrace{\qquad\qquad\qquad}_{\text{Bayes factor}} \underbrace{\qquad\qquad\qquad}_{\text{Prior odds ratio}}$$

Common approach (we may call it the “model selection approach”):

1. Take \hat{m} , e.g. $\hat{m} = \arg \max_m p(m | \mathbf{g})$
2. Use it as if it were “the true model”

This is an empirical Bayes approach

Alternative: model averaging (not focused here).

$$\hat{\mathbf{m}} = \arg \max_{\mathbf{m}} \{ \log p(\mathbf{g} | \mathbf{m}) + \log p(\mathbf{m}) \}$$

Main difficulty: computing the evidence

To simplify the notation: $\mathbf{f}_{(m)} \equiv (\mathbf{f}, \phi_{(m)}, \psi_{(m)})$

Common approach: a Laplace approximation (valid asymptotically)

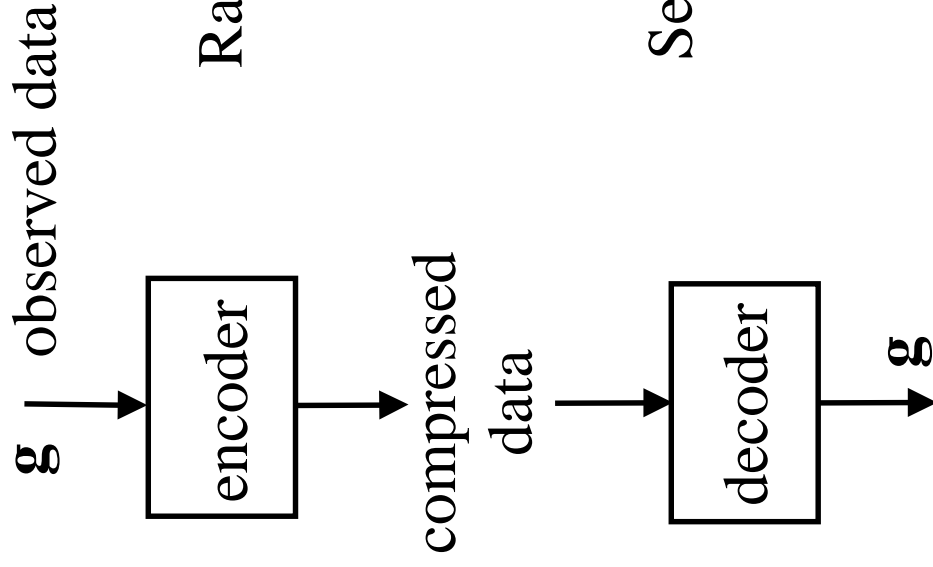
$$\log p(\mathbf{g} | \mathbf{m}) \approx \log p(\mathbf{g} | \hat{\mathbf{f}}_{(m)}^{\text{ML}}, \mathbf{m}) - \frac{\dim(\mathbf{f}_{(m)})}{2} \log n$$

Schwarz's Bayesian inference criterion (BIC)

Can be seen as an "order-penalized" log-likelihood.

25 The minimum description length (MDL) criterion

Introduction to MDL



Rationale: short code \Leftrightarrow good model

long code \Leftrightarrow bad model

code length \Leftrightarrow model adequacy

Several flavors: Rissanen 78, JRSS 87

Rissanen 96,

Wallace and Freeman, JRSS 87

26 Formalizing the MDL criterion

Scenario: we are given a set of models for the data

each model m is characterized by (unknown) “parameters” $\mathbf{f}_{(m)}$

$$\{p(\mathbf{g} | \mathbf{f}_{(m)}), m = m_1, m_2, \dots, m_K\}$$

no (or vague) prior information about $\mathbf{f}_{(m)}$

Goal: given data \mathbf{g} ,

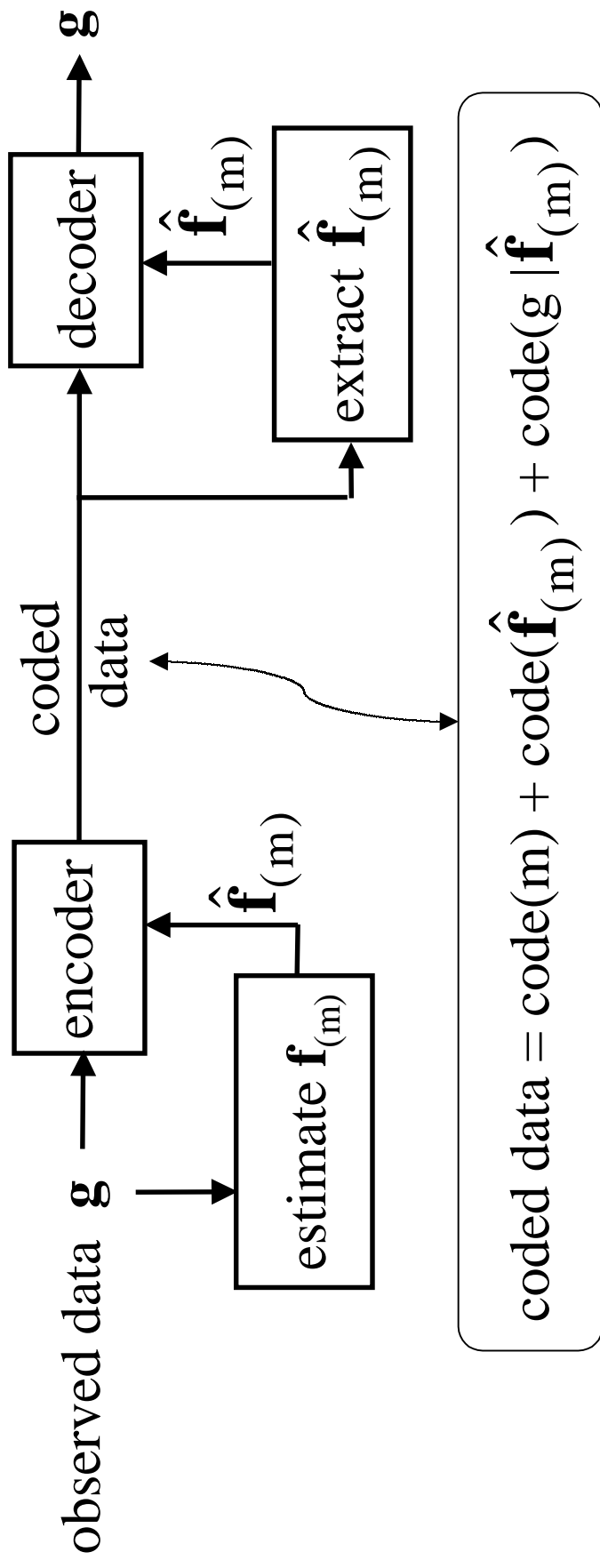
build the shortest possible code for \mathbf{g}

With $\mathbf{f}_{(m)}$ known, the shortest code-length for \mathbf{g} is (Shannon’s)

$$L(\mathbf{g} | \mathbf{f}_{(m)}) = -\log p(\mathbf{g} | \mathbf{f}_{(m)}, m)$$

However, $\mathbf{f}_{(m)}$ is, a priori, unknown; it has to be estimated.

Assumption: given $\mathbf{f}_{(m)}$, both encoder and decoder know how to build the same code



MDL principle: $\hat{\mathbf{f}}_{(m)}$ so that $\text{length}(\text{coded data})$ is shortest

28 Formalizing the MDL criterion

coded data = $\text{code}(m) + \text{code}(\hat{\mathbf{f}}_{(m)}) + \text{code}(\mathbf{g} | \hat{\mathbf{f}}_{(m)})$

$$L(m, \mathbf{f}_{(m)}, \mathbf{g}) = \underbrace{L(m)}_{\text{Usually constant}} + L(\mathbf{f}_{(m)} | m) + L(\mathbf{g} | \mathbf{f}_{(m)})$$

Usually constant

MDL criterion

$$\begin{aligned} (\hat{m}, \hat{\mathbf{f}}_{(\hat{m})})_{\text{MDL}} &= \arg \min_{m, \mathbf{f}_{(m)}} \{L(\mathbf{f}_{(m)}) + L(\mathbf{g} | \mathbf{f}_{(m)})\} \\ &= \arg \min_{m, \mathbf{f}_{(m)}} \{L(\mathbf{f}_{(m)}) - \log p(\mathbf{g} | \mathbf{f}_{(m)})\} \end{aligned}$$

Formalizing the MDL criterion

$$(\hat{m}, \hat{\mathbf{f}}_{(m)})_{\text{MDL}} = \arg \min_{m, \mathbf{f}_{(m)}} \{L(\mathbf{f}_{(m)}) - \log p(\mathbf{g} | \mathbf{f}_{(m)})\}$$

$L(\mathbf{f}_{(m)})$? Finite $L(\mathbf{f}_{(m)}) \Rightarrow$ truncate to finite precision: $\tilde{\mathbf{f}}_{(m)}$

High precision

$$-\log f(\mathbf{g} | \tilde{\mathbf{f}}_{(m)}) \approx -\log f(\mathbf{g} | \hat{\mathbf{f}}_{(m)}^{\text{ML}})$$

$$\text{but } L(\tilde{\mathbf{f}}_{(m)}) \nearrow$$

Low precision

$$-\log f(\mathbf{g} | \tilde{\mathbf{f}}_{(m)}) \text{ may be } \gg -\log p(\mathbf{g} | \hat{\mathbf{f}}_{(m)}^{\text{ML}})$$

$$L(\tilde{\mathbf{f}}_{(m)}) \searrow \text{ but}$$

Optimal compromise (under regularity conditions, and asymptotic)

$$L(\text{each component of } \mathbf{f}_{(m)}) = \frac{1}{2} \log(n)$$

n is the sample size
from with the parameter
is estimated
(growth rate of Fisher info.)

30 Formalizing the MDL criterion

Under regularity conditions

$$(\hat{m}, \hat{\mathbf{f}}_{(\hat{m})})_{\text{MDL}} = \arg \min_{m, \mathbf{f}_{(m)}} \left\{ -\log p(\mathbf{g} | \mathbf{f}_{(m)}) + \frac{\dim(\mathbf{f}_{(m)})}{2} \log(n) \right\}$$

Coincides with Schwarz's BIC.

the dimension of $\mathbf{f}_{(m)}$



This is the “classical” (1978) MDL criterion.

It is not valid in our problems.

Basically: each “parameter” is not estimated from all the “data”.

Quick and dirty engineering solution:

use the “natural” description length

Back to the image restoration problem

Likelihood function $p(\mathbf{g} | \mathbf{f}, \sigma^2) = \mathcal{N}(\mathbf{g} | \mathbf{f}, \sigma^2 \mathbf{I})$

Prior $p(\mathbf{f} | \Psi_{(m)}, m) = \mathcal{N}(\mathbf{f} | \mathbf{0}, [\mathbf{A}(\Psi_{(m)})]^{-1})$

$p(\mathbf{g}, \mathbf{f} | \sigma^2, \Psi_{(m)}, m) \propto p(\mathbf{g} | \mathbf{f}, \sigma^2) p(\mathbf{f} | \Psi_{(m)}, m)$

$$(\hat{m}, \hat{\Psi}_{(\hat{m})}, \hat{\sigma}^2)_{\text{MDL}} = \arg \min_{m, \Psi_{(m)}} \{-\log p(\mathbf{g}, \mathbf{f} | \sigma^2, \Psi_{(m)}, m) + L(\Psi_{(m)})\}$$

Can be seen as an MDL criterion with missing data (\mathbf{f} is unknown).
Implemented via an EM-type algorithm

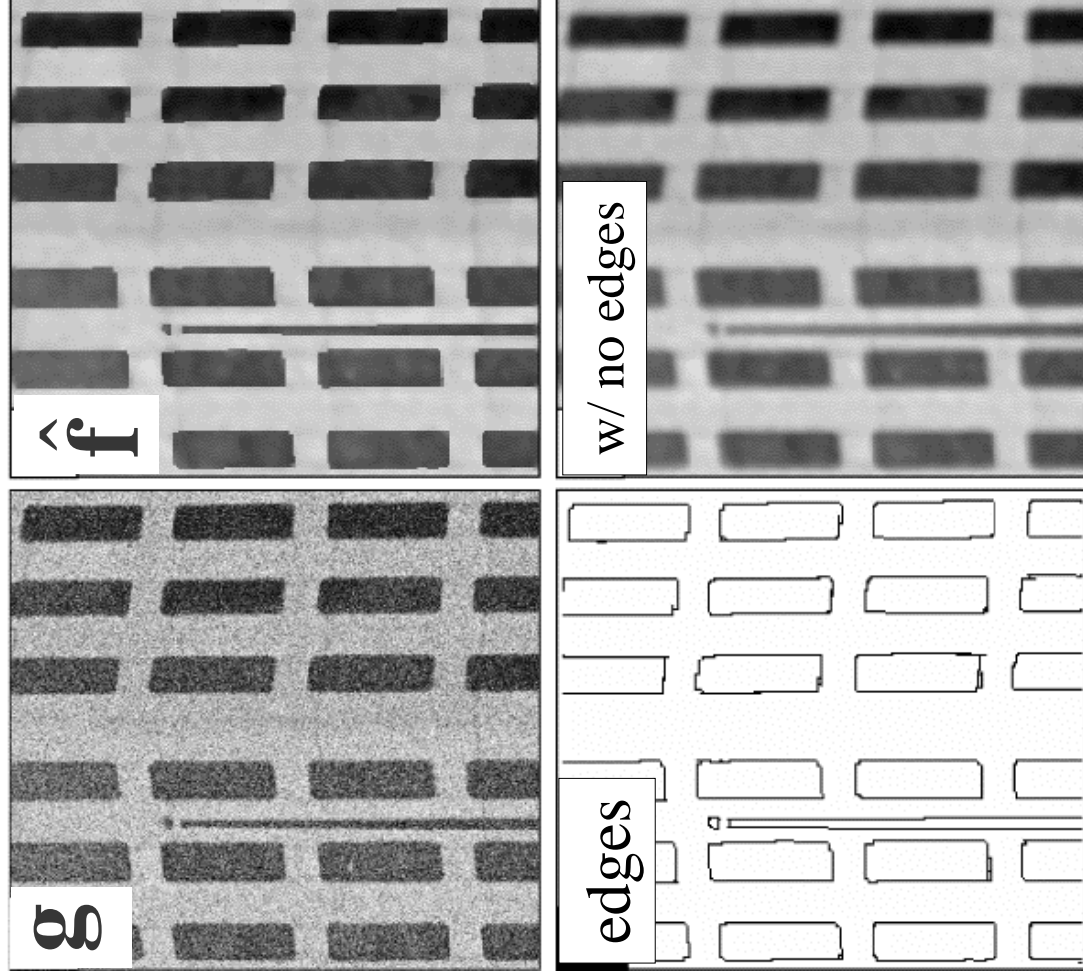
Given $\hat{\sigma}^2$ and $\hat{\Psi}_{(m)}$

$$\hat{\mathbf{f}} = [\hat{\sigma}^2 \mathbf{A}(\hat{\Psi}_{(m)}) + \mathbf{I}]^{-1} \mathbf{g}$$

$$L(\Psi_{(m)}) = m \log(N)$$

Natural code-length for m edge locations on an N -pixels image

32 Image restoration examples



Totally unsupervised

μ , σ^2 , number of edges,
their locations

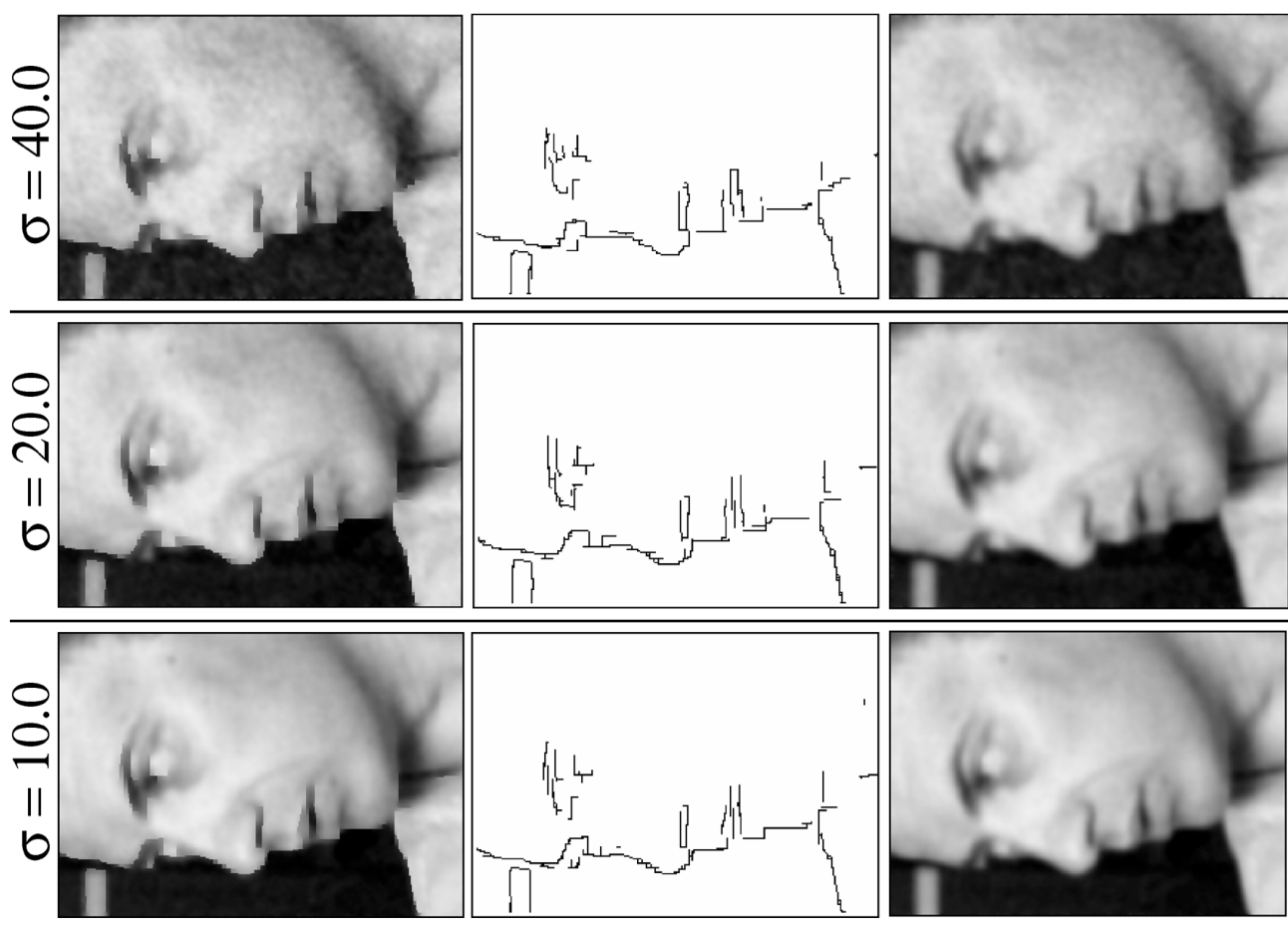
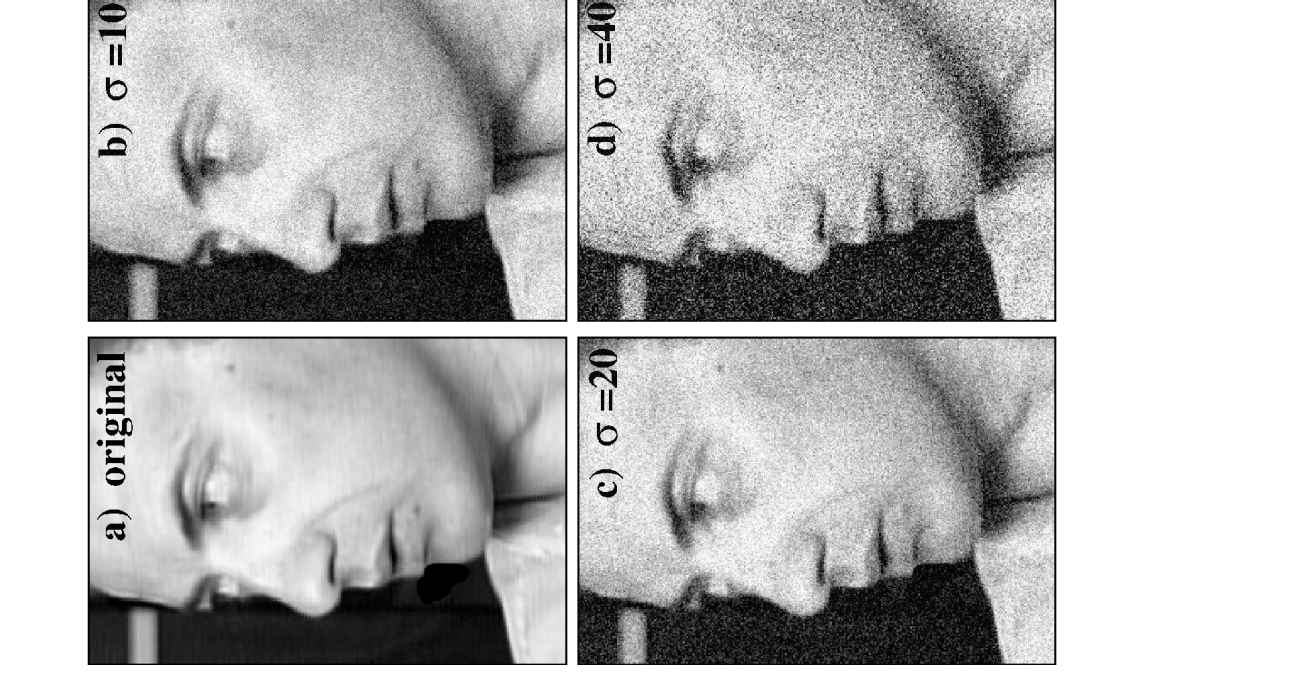


all unknown

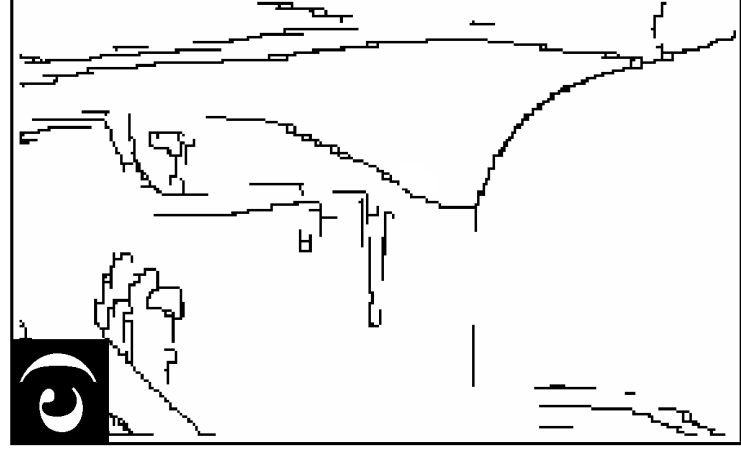
Implemented via an EM-type algorithm, together with a continuation method, and a pseudo-likelihood approximation

M. Figueiredo and J. Leitão, “Unsupervised image restoration and Edge location using Gauss-Markov random fields and the MDL principle”, *IEEE Transactions on Image Processing*, vol. 6, pp. 1089-1102, 1997.

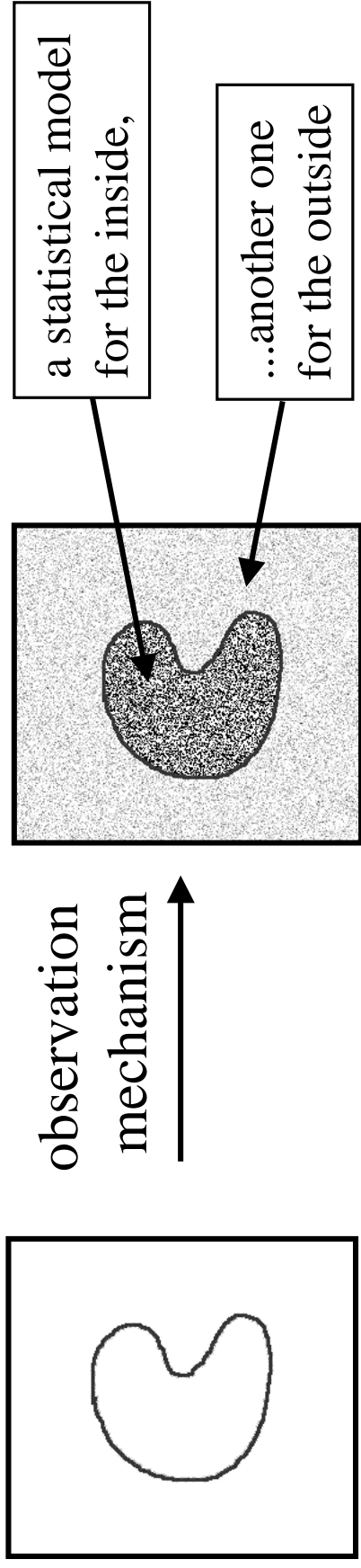
Further image restoration examples



34 Further image restoration examples



35 Back to the contour estimation problem



\mathbf{f} , contour description

\mathbf{g} , observed image

$\mathbf{f}_{(m)} = \mathbf{B}_{(m)} \boldsymbol{\theta}_{(m)} \longrightarrow$ a 2-D spline curve with m control points

$$(\hat{m}, \hat{\boldsymbol{\theta}}_{(\hat{m})})_{MDL} = \arg \min_{m, \boldsymbol{\theta}_{(m)}} \{ -\log p(\mathbf{g} | \boldsymbol{\theta}_{(m)}) + L(\boldsymbol{\theta}_{(m)}) \}$$

$$p(\mathbf{g} | \boldsymbol{\theta}_{(m)}, \boldsymbol{\phi}_{in}, \boldsymbol{\phi}_{out}) = \prod_{i \in \text{inside}(\mathbf{f}_{(m)})} p(\mathbf{g}_i | \boldsymbol{\phi}_{in}) \prod_{i \in \text{outside}(\mathbf{f}_{(m)})} p(\mathbf{g}_i | \boldsymbol{\phi}_{out})$$

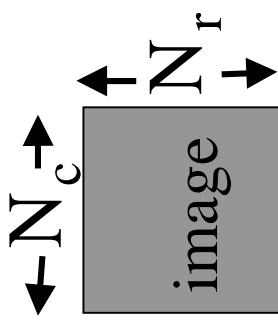
36 Contour estimation problem

Regarding $\phi_{\text{in}}, \phi_{\text{out}}$ → minimize also w.r.t. them
→ adopt a (vague, conjugate) prior, and integrate them out

What about $L(\theta_{(m)})$?

Natural code length for m control points

$$L(\theta_{(m)}) = m[\log(N_r) + \log(N_c)] = m \log(N)$$



$\|B_{(m)}\|_{\infty} = 1$ → because of the partition of unity of B-splines
quantizing control points at pixel resolution
guarantees also pixel resolution for the contour

37 Contour estimation problem

For each m , $\hat{\Theta}_{(m)}^{\text{ML}}$ is obtained by a gradient-projection-type algorithm

Then,

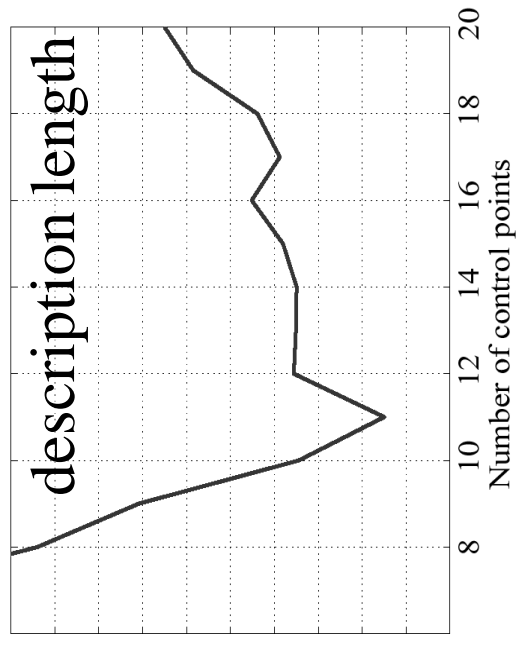
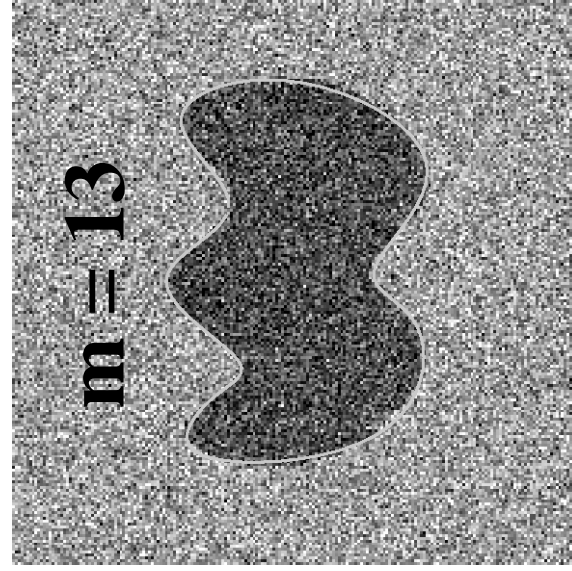
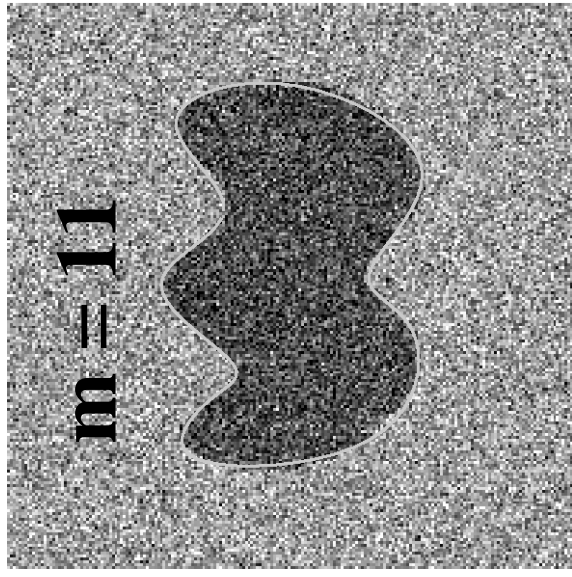
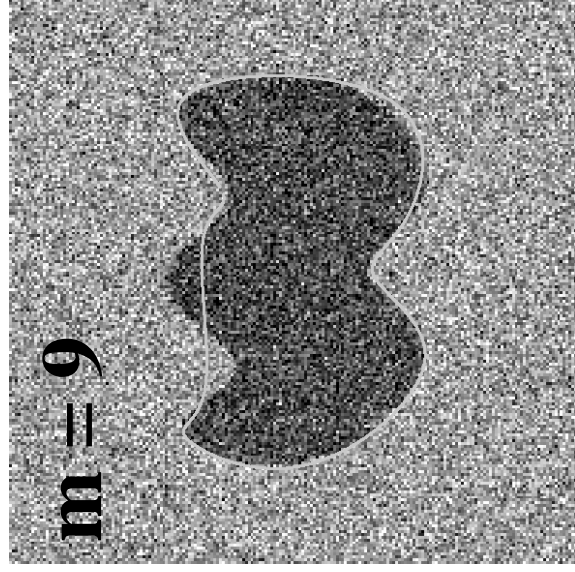
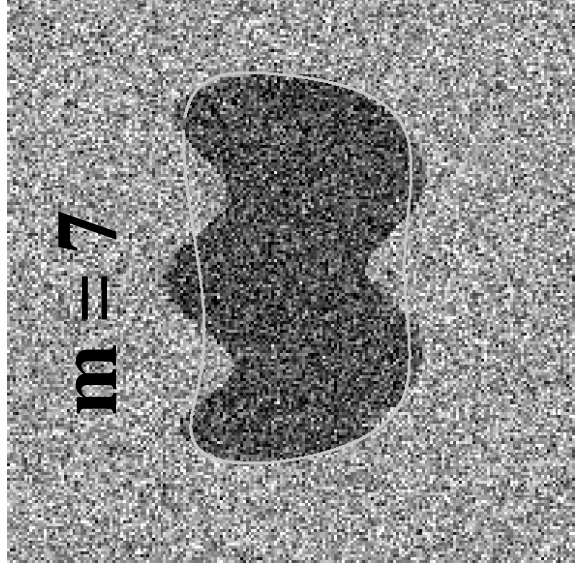
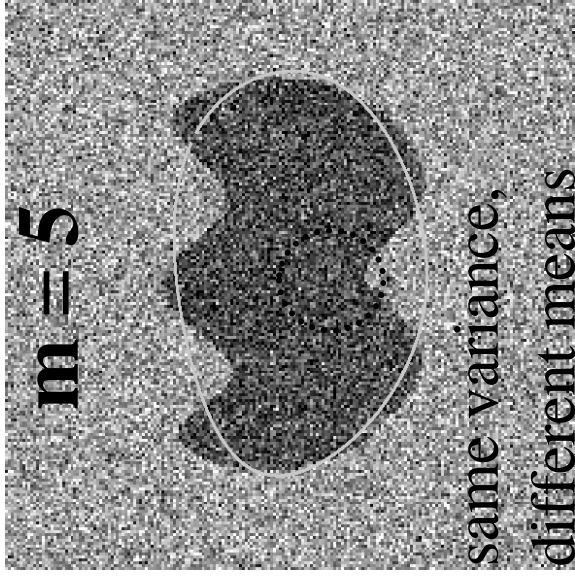
$$\hat{m}_{\text{MDL}} = \arg \min_m \{ -\log p(\mathbf{g} | \hat{\Theta}_{(m)}^{\text{ML}}) + m \log(N) \}$$

Examples:

- synthetic images: Gaussian regions of different mean and variance
- ultrasound images: Rayleigh regions of different variance
- NMR: Riccian regions of different mean and variance

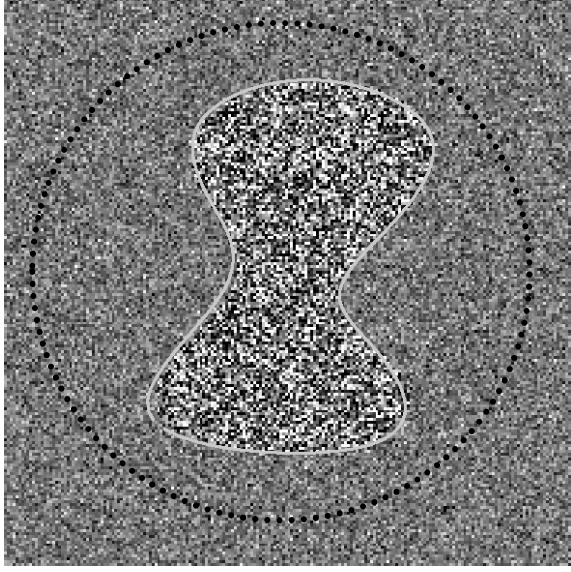
See: M. Figueiredo, J. Leitão, and A. Jain, “Adaptive B-splines and boundary estimation”, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR-97), pp. 724-729, San Juan, Puerto Rico, 1997.

38 Contour estimation examples

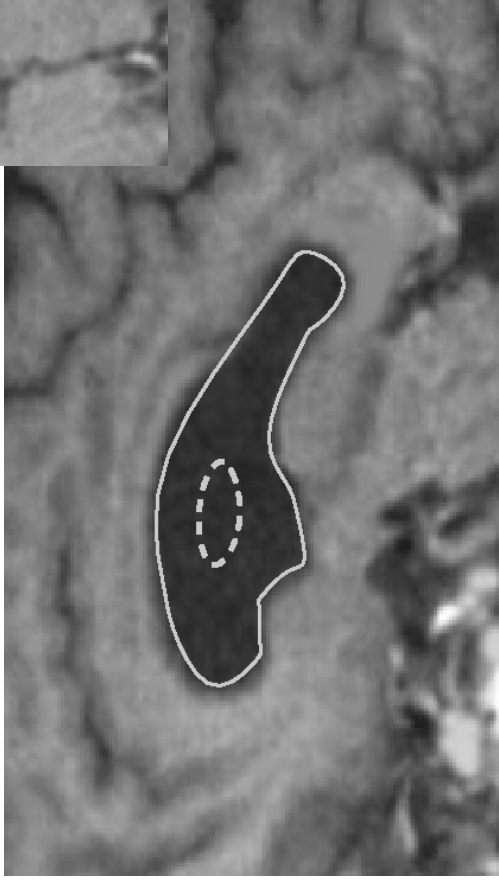
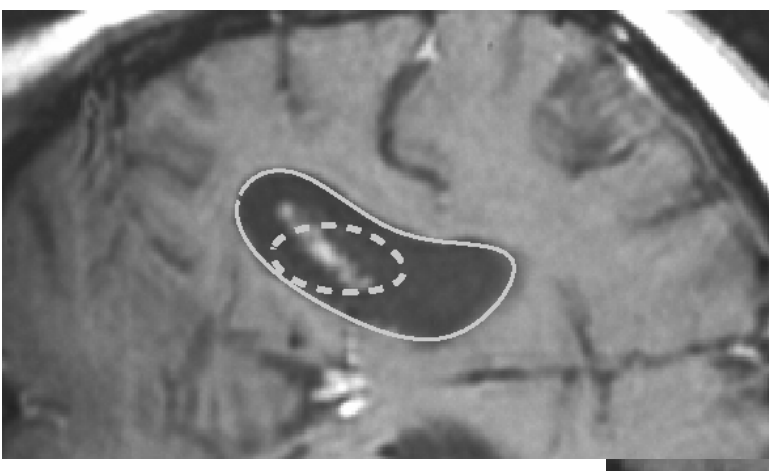
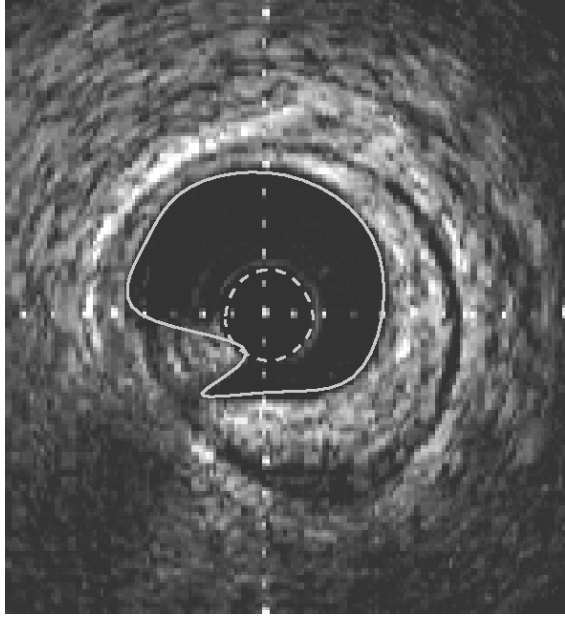
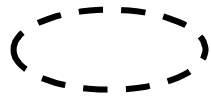


39 Contour estimation examples

same mean,
different variances



initializations



40 The incompleteness issue in MDL

Classical MDL (Rissanen, 1978)

$$(\hat{m}, \hat{\mathbf{f}}_{(\hat{m})})_{\text{MDL}} = \arg \min_{m, \mathbf{f}_{(m)}} \left\{ \frac{m}{2} \log(n) - \log p(\mathbf{g} | \mathbf{f}_{(m)}) \right\}$$

Problems: not valid for small samples (asymptotic criterion);

assumes a redundant (called *incomplete*) code (Rissanen, 1996).

Incompleteness: in code-length $-\log p(\mathbf{g} | \mathbf{f}_{(m)})$

1. Transmitter estimates $\mathbf{f}_{(m)}$ from \mathbf{g}
2. Transmitter sends $\hat{\mathbf{f}}_{(m)}$ to receiver
3. Now, receiver knows that the only possible \mathbf{g} 's are those that can lead to that $\hat{\mathbf{f}}_{(m)}$
Code only needed for those \mathbf{g} 's

The incompleteness issue in MDL

Example: $\mathbf{g} = \{g_1, g_2, \dots, g_{10}\}$ $g_i \in \{0, 1\} \implies \mathbf{g} \in \{0, 1\}^{10}$

10 i.i.d. observations of a Bernoulli source

$\text{Prob}(g_i=1) = f$

$$1. \text{ Transmitter computes } \hat{f} = \frac{1}{10} \sum_{i=1}^{10} g_i$$

2. Transmitter sends \hat{f} to receiver (rational number)

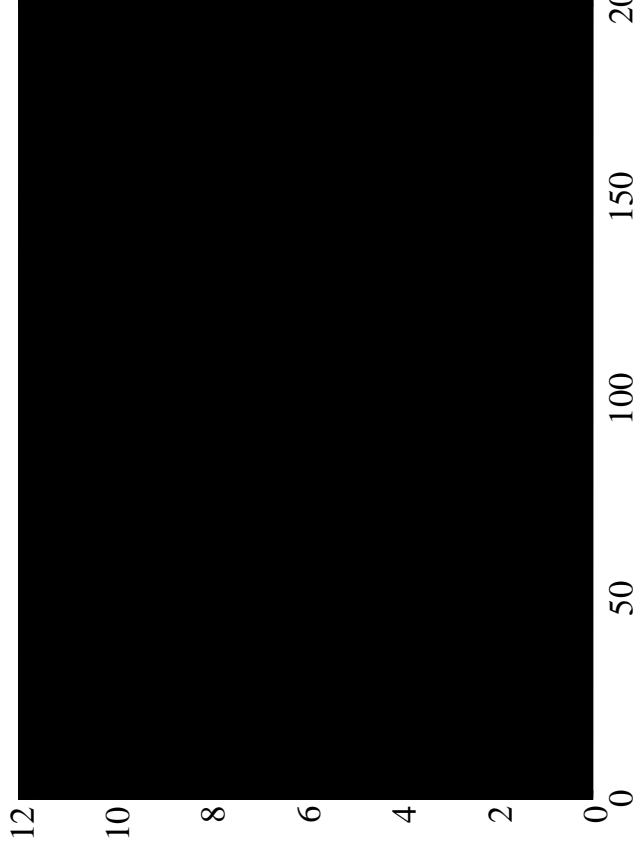
3. Not all \mathbf{g} 's need be coded.

Only those that have 10 \hat{f} ones: $L(\mathbf{g}) = \log \binom{10}{10\hat{f}}$

Concrete example: $\mathbf{g} = \{1, 0, 0, 0, 1, 0, 1, 0, 0, 0\}$, 3 ones $\longrightarrow \hat{f} = \frac{3}{10}$

$$\log_2 \binom{10}{3} = 6.9 \text{ bits} \quad \text{versus} \quad -\log_2 p(\mathbf{g}|\hat{f}) = -\log_2 [0.3^3 (1-0.3)^7] \approx 8.8 \text{ bits}$$

Basic question: given a sequence of independent Poisson counts



$$\mathbf{g} = \{g_1, \dots, g_N\}$$

Is there a change in the parameter ?
If so, where ?

$$m_0 : p(\mathbf{g} | \lambda) = \prod_{i=1}^N \text{Poisson}(g_i | \lambda)$$

We have $N+1$ models:

$$m_1, \dots, m_N : p(\mathbf{g} | \lambda_a, \lambda_b, m_k) = \prod_{i=1}^k \text{Poisson}(g_i | \lambda_a) \prod_{i=k+1}^N \text{Poisson}(g_i | \lambda_b)$$

An MDL criterion to segment Poisson data

Fact: g_1, \dots, g_N be N i.i.d. Poisson counts (arbitrary parameter).

Let $t = \sum g_i$ be the total count

Conditioned on t , the g_i 's follow a multinomial distribution

$$p(g_1, \dots, g_N | t) = \text{Multi}(g_1, \dots, g_N | t; \frac{1}{N}, \dots, \frac{1}{N})$$

“each g_i is expected to get a $1/N$ fraction of the total count t ”

What if $g_1, \dots, g_i \sim$ i.i.d. $\text{Poisson}(\lambda_a)$ $g_{i+1}, \dots, g_N \sim$ i.i.d. $\text{Poisson}(\lambda_b)$

$$p(g_1, \dots, g_N | t) = \text{Multi}(g_1, \dots, g_N | t; \underbrace{\frac{2\rho}{N}, \dots, \frac{2\rho}{N}}_{1 \text{ to } k}, \underbrace{\frac{2(1-\rho)}{N}, \dots, \frac{2(1-\rho)}{N}}_{k+1 \text{ to } N})$$

$$\rho = \frac{\lambda_a}{\lambda_a + \lambda_b}$$

Approach: rather than base the test on the Poisson model, condition on the total count, and use the multinomial

The $N+1$ models:

$$m_0 : p(\mathbf{g}_1, \dots, \mathbf{g}_N | \mathbf{t}) = \binom{\mathbf{t}}{\mathbf{g}_1 \mathbf{g}_1 \dots \mathbf{g}_N} \left(\frac{1}{N} \right)^{\mathbf{t}}$$

$$m_k : p(\mathbf{g}_1, \dots, \mathbf{g}_N | \mathbf{t}) = \binom{\mathbf{t}}{\mathbf{g}_1 \mathbf{g}_1 \dots \mathbf{g}_N} \left(\frac{2\rho_k}{N} \right)^{\sum_{i=1}^k \mathbf{g}_i} \left(\frac{2(1-\rho_k)}{N} \right)^{\sum_{i=k+1}^N \mathbf{g}_i}$$

Note that m_0 is nested into all the m_k 's

$$\lambda_a = \lambda_b \Rightarrow \rho = \frac{\lambda_a}{\lambda_a + \lambda_b} = \frac{1}{2} \Rightarrow \frac{2\rho}{N} = \frac{2(1-\rho)}{N} = \frac{1}{N}$$

Our MDL approach:

1. given $\mathbf{g} = \{g_1, \dots, g_N\}$ compute $t = \sum g_i$
2. transmit t (e.g., with an arbitrary technique, e.g. Elias code)
3. compute the description length under model m_0

$$m_0 : L(\mathbf{g} | t, m_0) = -\log \binom{t}{g_1 g_1 \dots g_N} + t \log N$$

This is the total description length, because there is no parameter

4. compute the description length under models m_1, \dots, m_N

4. compute the description length under models m_1, \dots, m_N

Under m_k : $-\log(\text{Multinomial with two parameters: } \frac{2\rho_k}{N} \text{ and } \frac{2(1-\rho_k)}{N})$

First thing: estimate ρ_k : $\hat{\rho}_k = \frac{1}{t} \sum_{i=1}^k g_i = \frac{s_k}{t}$

Code length for $\hat{\rho}_k$?

Since t was already sent, we only need to encode s_k

$s_k \in \{0, 1, 2, \dots, t\} \longrightarrow L(\rho_k) = L(s_k) = \log(t+1)$

Should we now plug $\hat{\rho}_k$ in the multinomial above ?

Attention ! Incompleteness. Receiver already knows s_k

An MDL criterion to segment Poisson data

Best code for $\mathbf{g} = \{g_1, \dots, g_N\}$

...knowing that $\sum_{i=1}^k g_i = s_k$..and consequently that $\sum_{i=k+1}^N g_i = t - s_k$

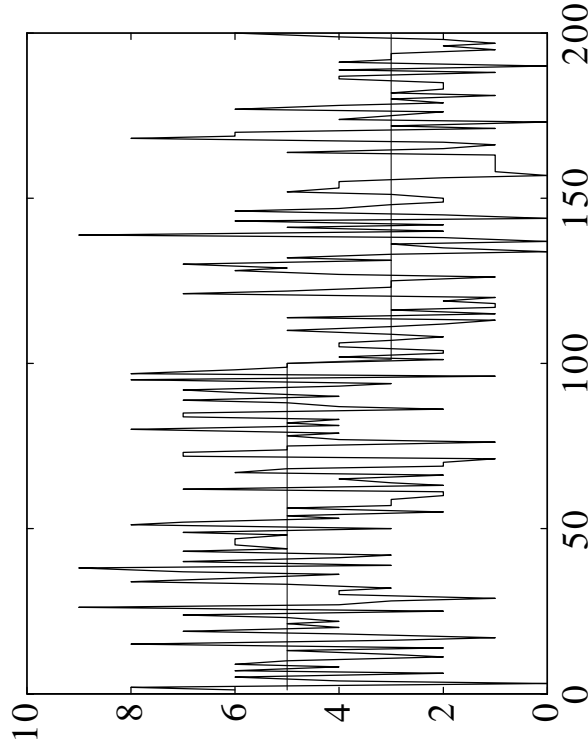
$$m_k : L(\mathbf{g} | t, m_k, s_k) = -\log \text{Multi}(g_1, \dots, g_k | s_k, \frac{1}{k}, \dots, \frac{1}{k}) \\ -\log \text{Multi}(g_{k+1}, \dots, g_N | (t - s_k), \frac{1}{N-k}, \dots, \frac{1}{N-k})$$

Finally, we choose the shortest of the $N+1$ total description lengths

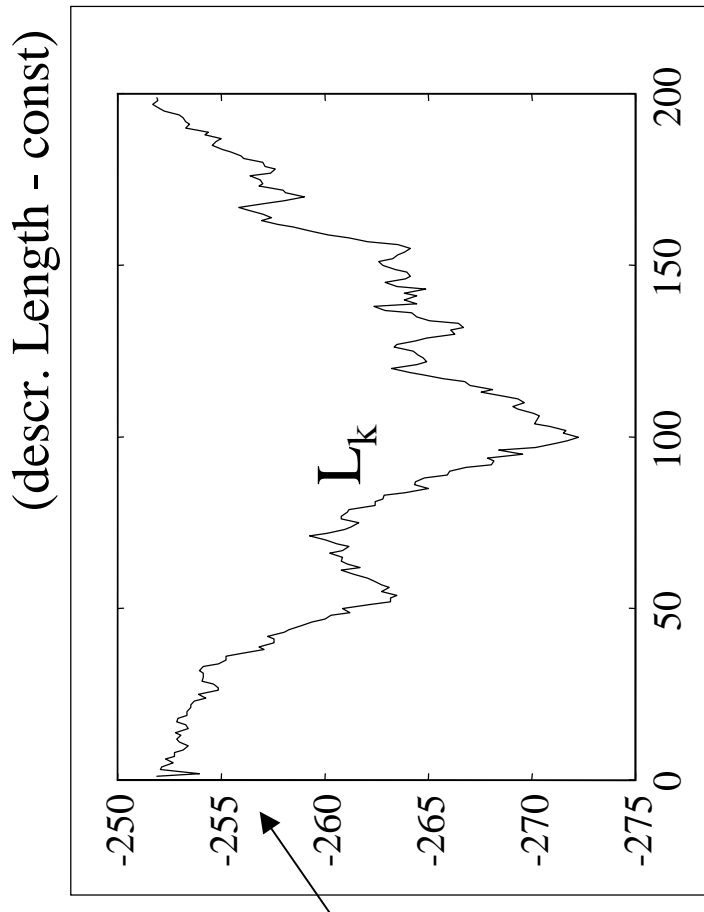
$$L_0 = L(\mathbf{g} | t, m_0) \quad L_k = \log(t+1) + L(\mathbf{g} | t, m_k, s_k), \quad k = 1, \dots, N$$

Surprisingly (or maybe not): this MDL criterion coincides with the Bayes model selection criterion with a uniform prior on $\rho \in [0,1]$

48 An MDL criterion to segment Poisson data



Minimum at $k = 100$
the correct value (in this case)



$$L_0 = -256$$

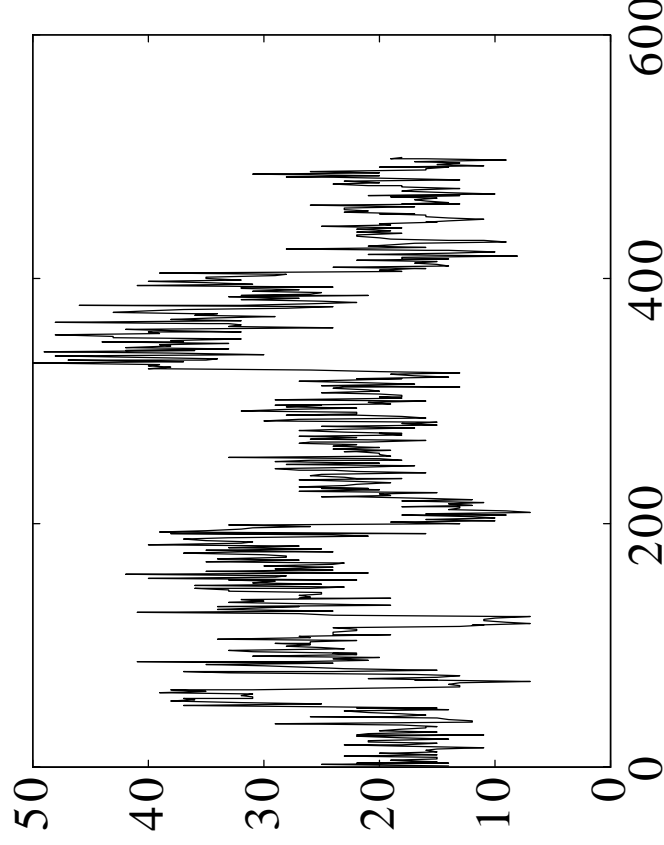
49 An MDL criterion to segment Poisson data

What about multiple change points ?

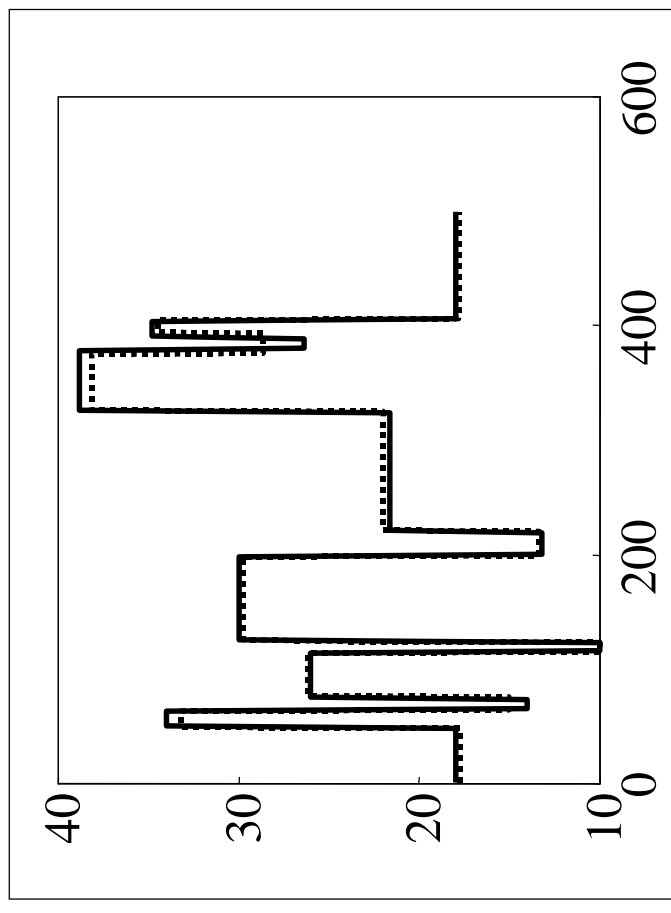
Simply take each segment and apply the criterion again (recursively)

Example:

Observed counts



— true intensity function
..... estimate



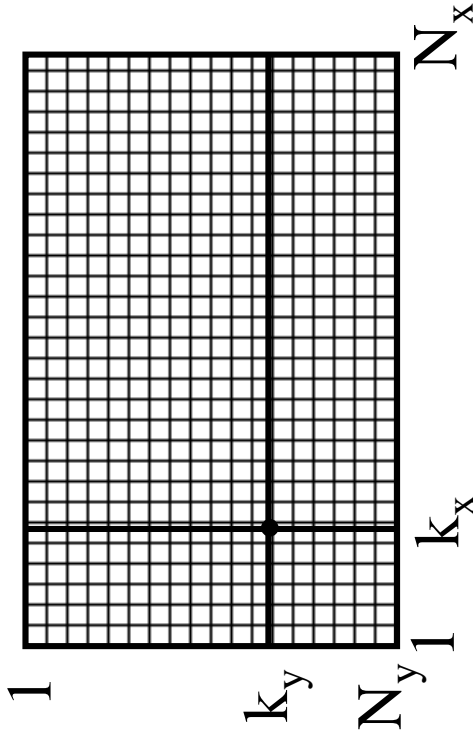
See: R. Nowak and M. Figueiredo, "Unsupervised progressive parsing of Poisson fields using minimum description length criteria", in Proceedings of the IEEE International Conference on Image Processing (ICIP-99), Kobe, Japan, October 1999 (to appear).

50 Segmenting Poisson images

Basic building block:

In 1D, we looked for the best (if any) change point

In 2D, we look for the best (if any) segmentation into two/four rectangles of the following type:



Same multinomial-based criterion:

- condition on the total count
- number of competing models:
 $(1 + (N_x + 1)(N_y + 1))$
no segmentation,
all possible 4-segmentations,
and all possible 2-segmentations

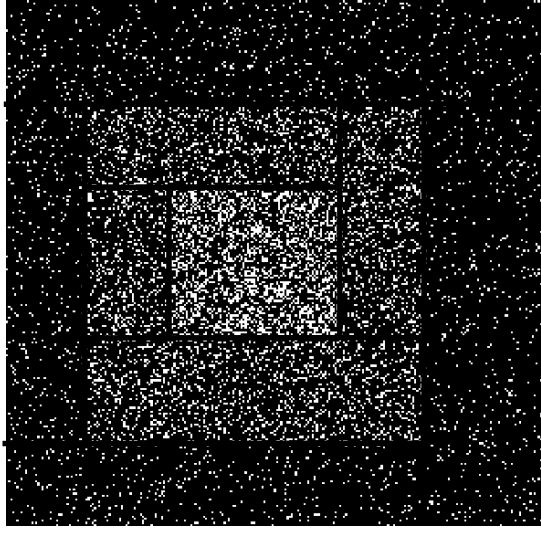
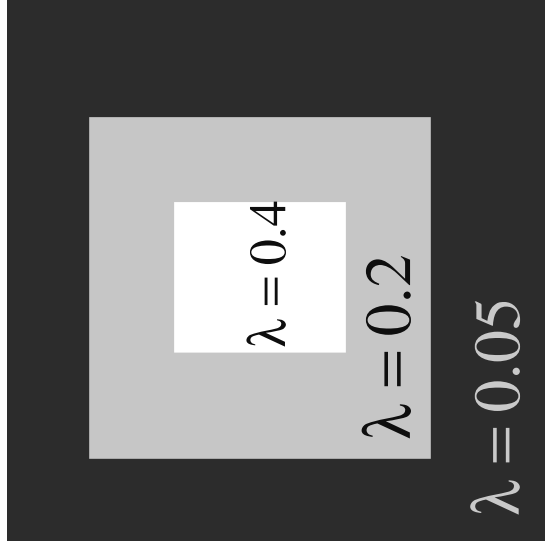
To fully segment an image: apply criterion recursively.

Segmenting Poisson images: an example

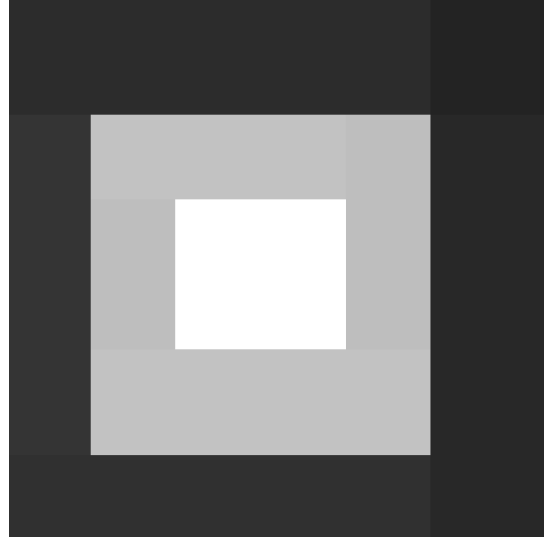
Example:

250*250 pixels/bins

counts



segmentation

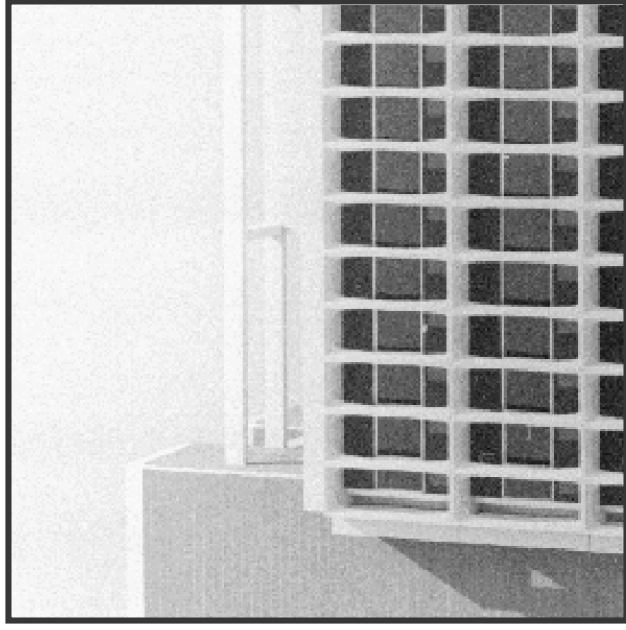


intensity estimate $\hat{\lambda}$

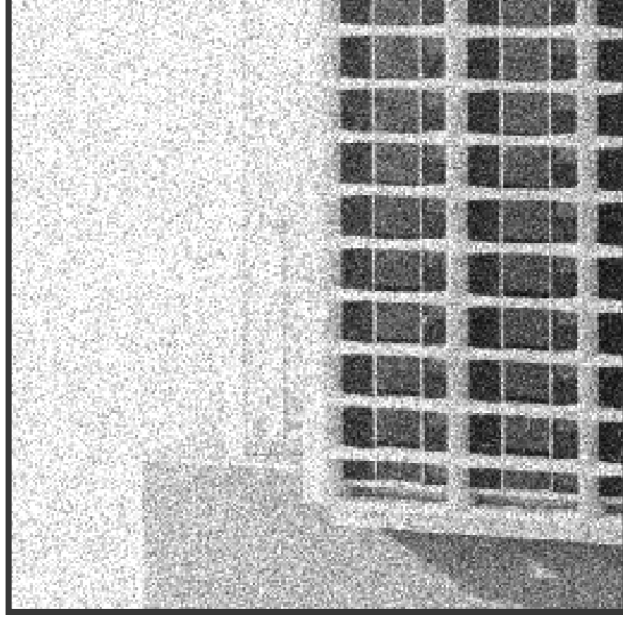
See: R. Nowak and M. Figueiredo, "Unsupervised progressive parsing of Poisson fields using minimum description length criteria", in Proceedings of the IEEE International Conference on Image Processing (ICIP-99), Kobe, Japan, October 1999 (to appear).

52 Segmenting Poisson images: another example

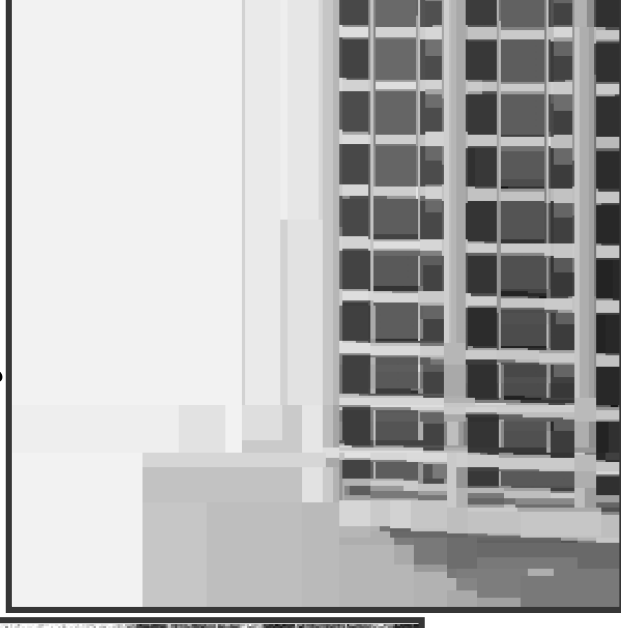
True intensity λ



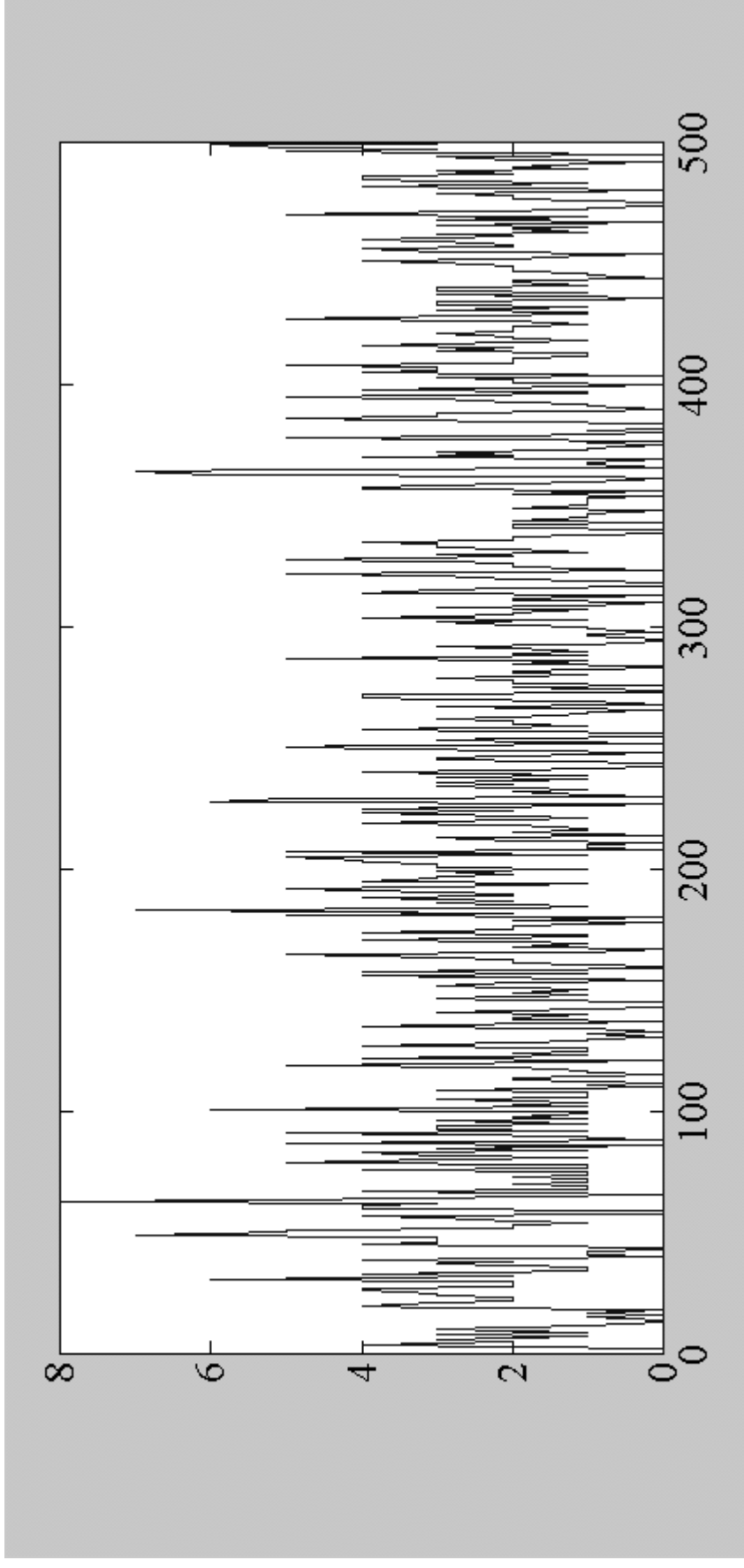
counts



intensity estimate $\hat{\lambda}$

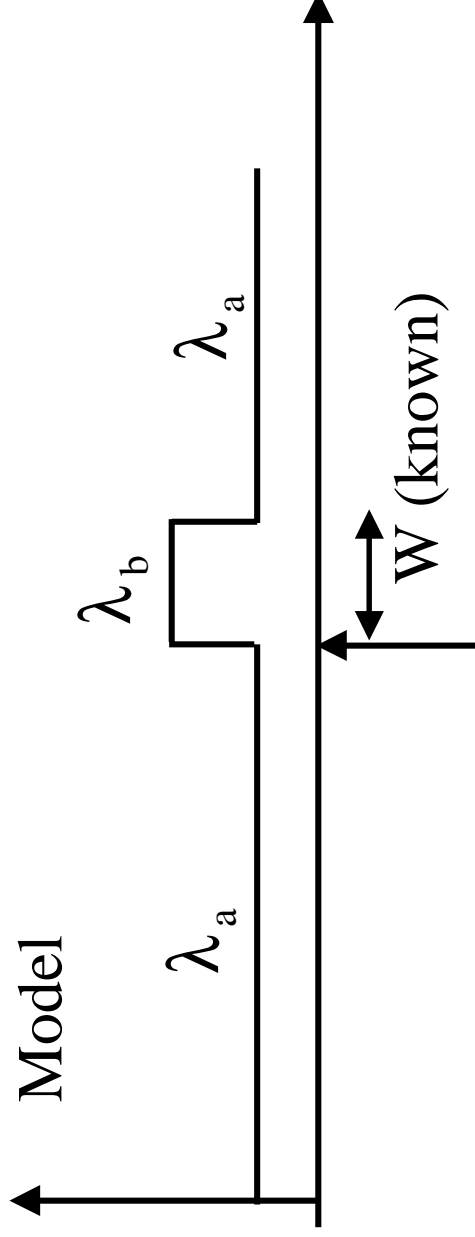


53 Singularities in Poisson data (“bursts”)



Question: given a sequence of Poisson observations,
...is there a “small” region of higher intensity than the background ?

54 Singularities in Poisson data (“bursts”)



Unknown position

λ_a, λ_b unknown (of course)

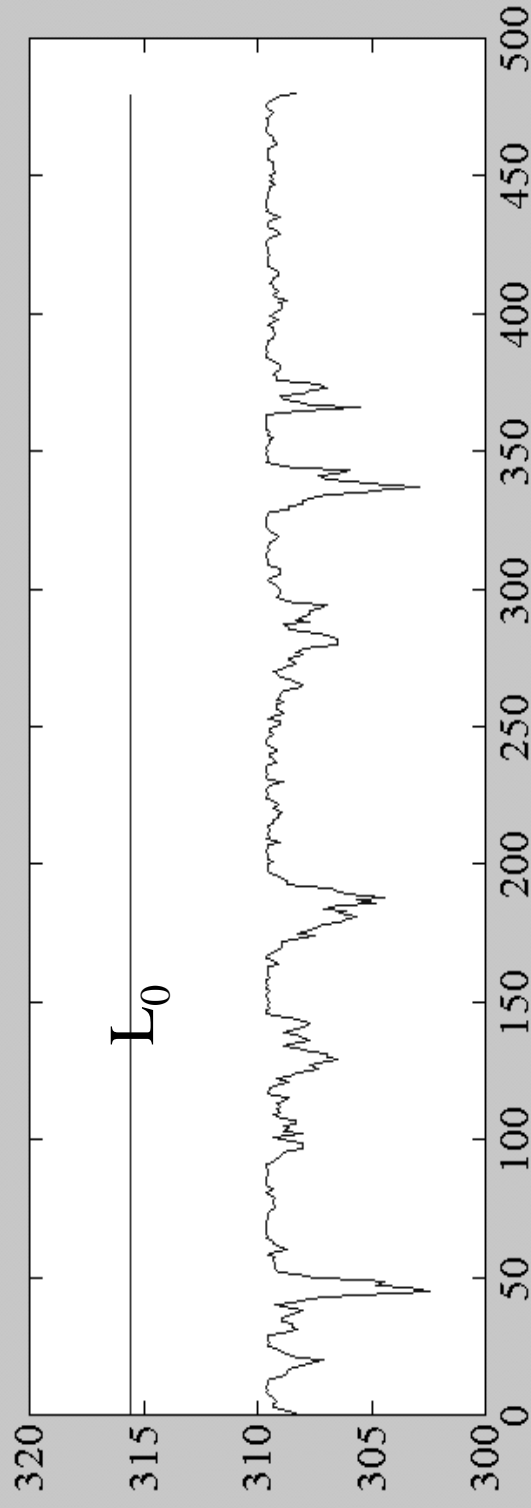
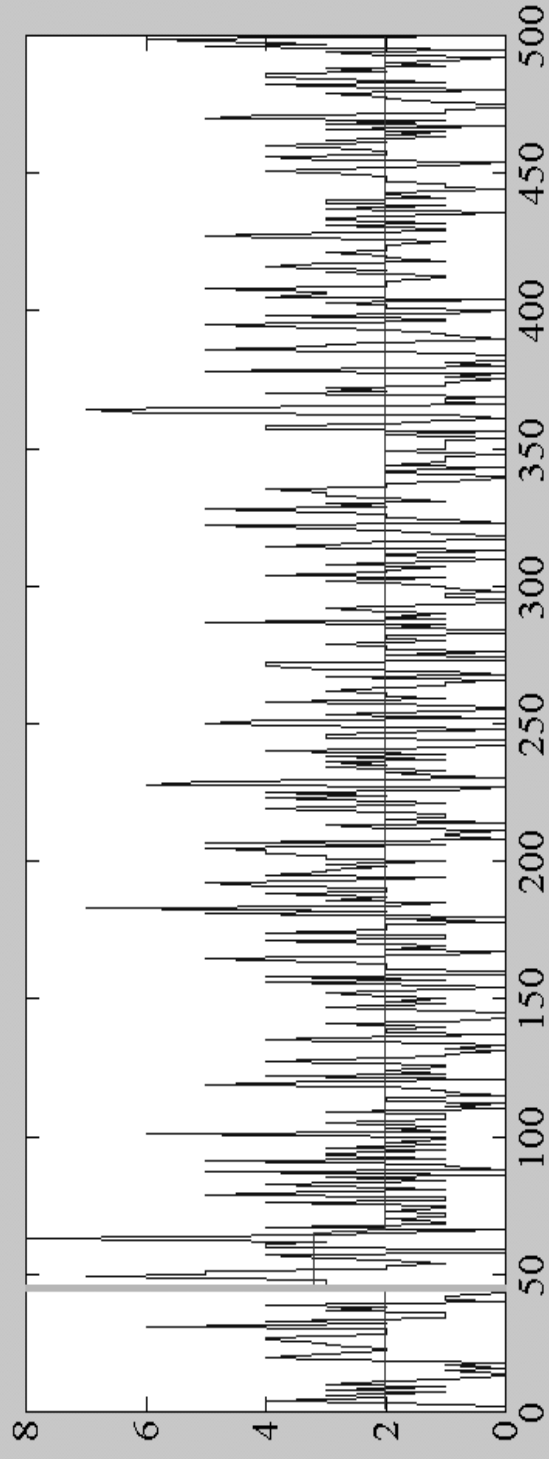
Approach as above: condition on the total count,
work with the multinomials

Again, equivalent to Bayes criterion with flat prior on $\rho = \frac{\lambda_a}{\lambda_a + \lambda_b}$

55 Singularities in Poisson data (“bursts”): example

$$\lambda_a = 2$$

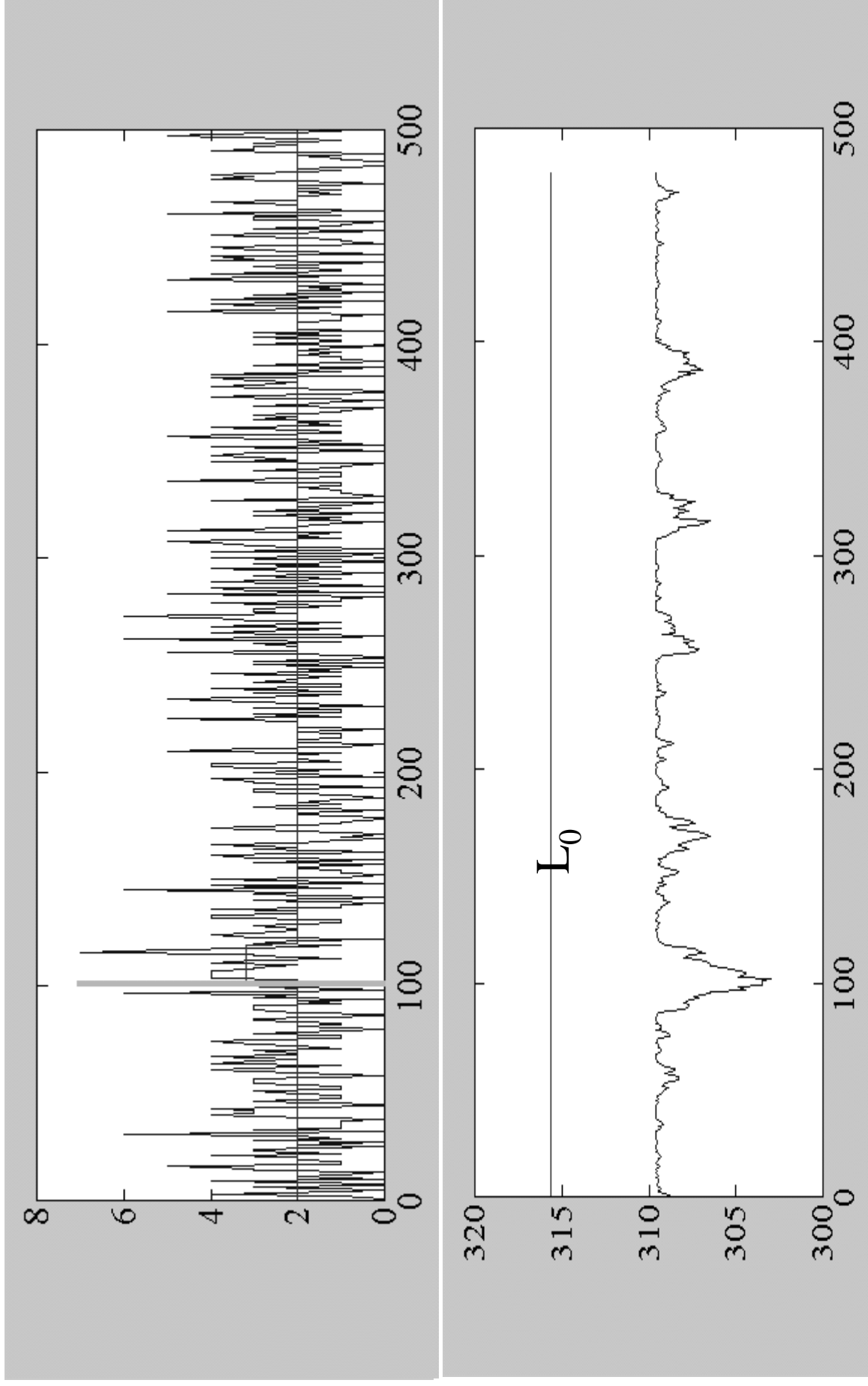
$$\lambda_b = 3$$



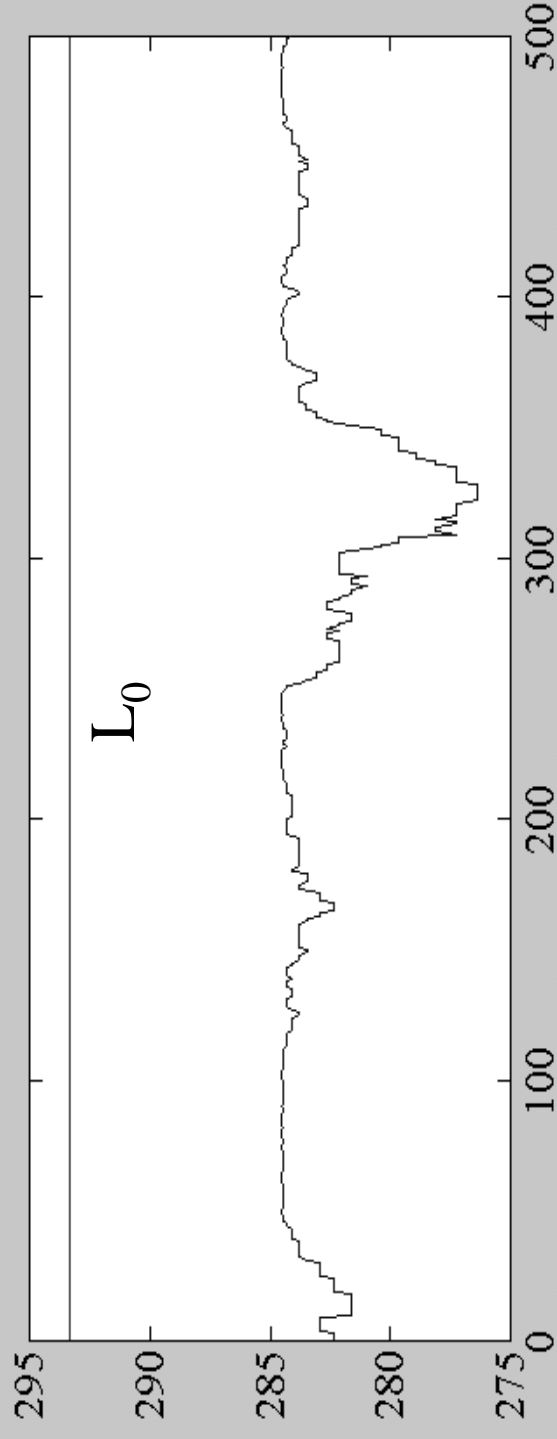
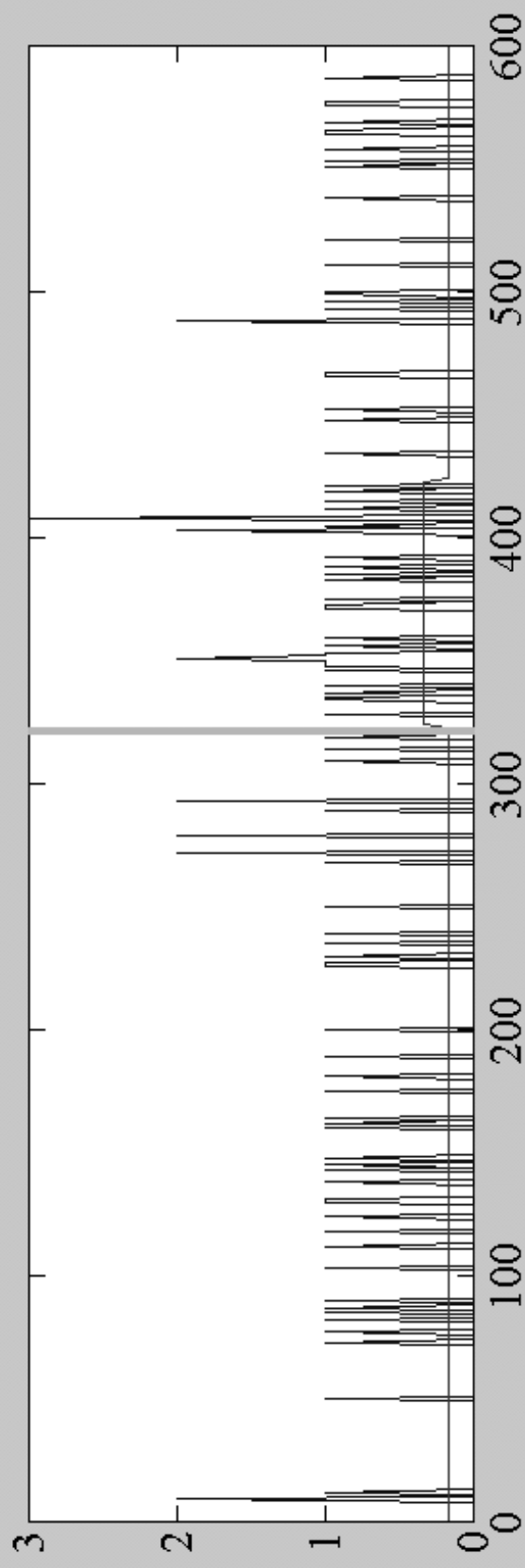
56 Singularities in Poisson data (“bursts”): example

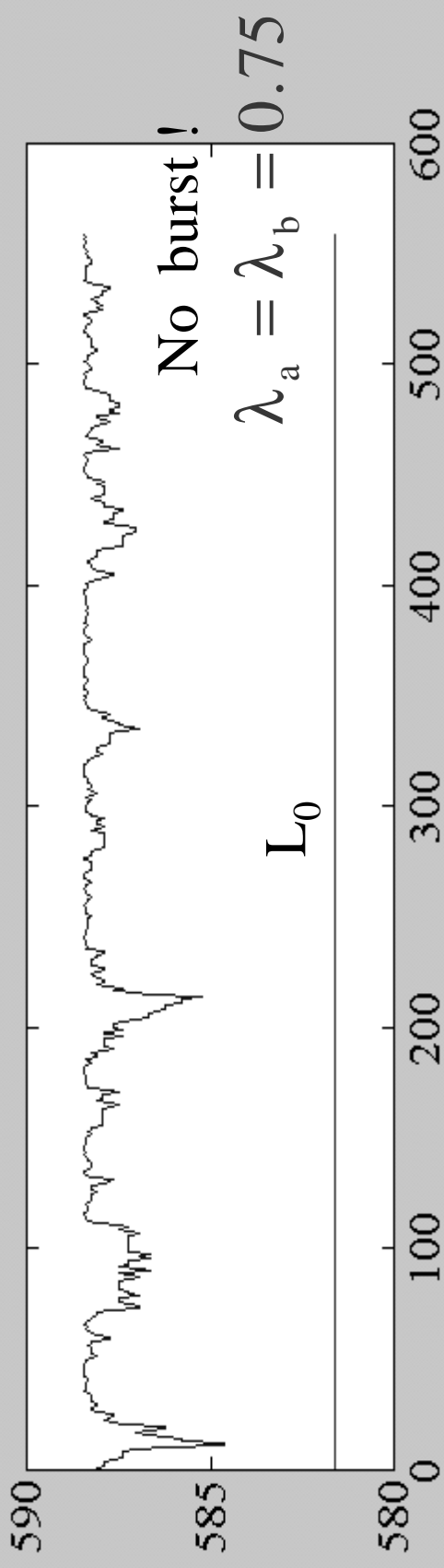
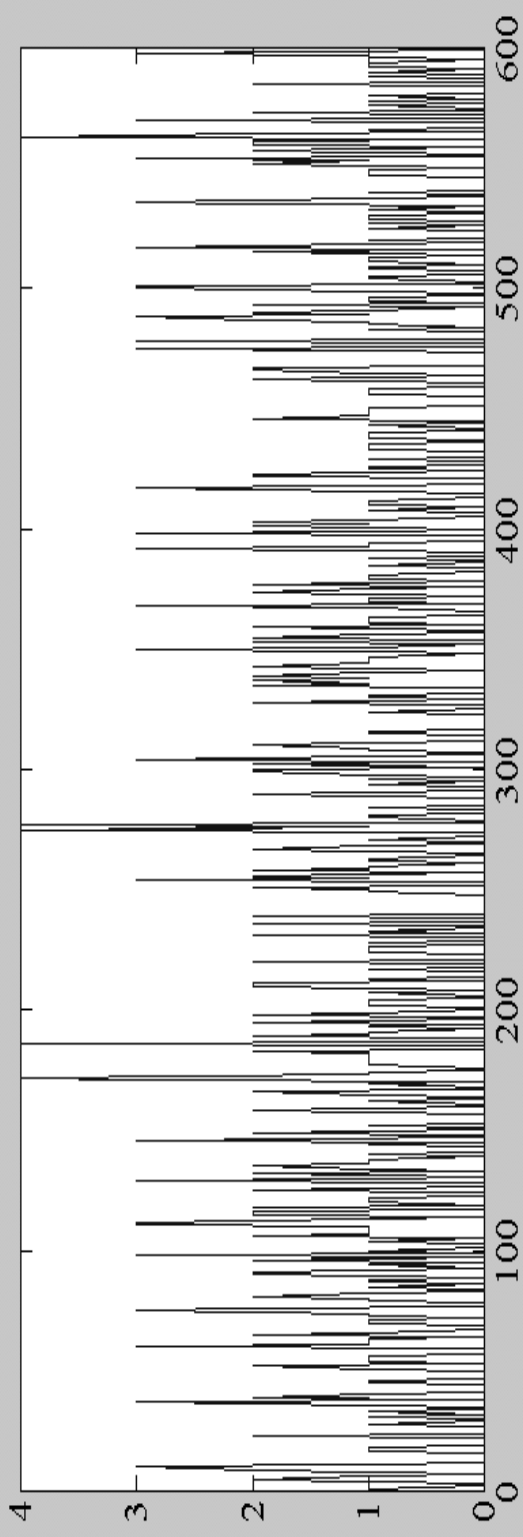
$$\lambda_a = 2$$

$$\lambda_b = 3$$



57 Singularities in Poisson data (“bursts”): example, low counts





- Most image analysis problems can/should be addressed as statistical (namely Bayesian) inference.
- Many such problems involve model selection issues. Available approaches:
 - Bayesian model selection
 - Minimum description length (MDL) model selection
 - Cross-validation / bootstrapping methods are computationally too heavy
- We addressed 3 examples of model selection problems in image analysis using MDL: image restoration, contour estimation, Poisson image segmentation
- In many applications, a natural code (rather than an optimal one) yields excellent results.
- If correctly formulated, MDL coincides with Bayesian model selection using reference priors.