Multiscale Data Analysis through Optimal Segmentation

Jeffrey.D.Scargle@nasa.gov Space Science Division NASA Ames Research Center

Brad Jackson, Mathematics Department, San Jose State University Jay Norris, NASA Goddard Spaceflight Center Michael Way, Paul Gazis, Chris Henze, Creon Levit, NASA-ARC Adam Roy, San Francisco State University Mahmoud Quweidar, University of Texas, Brownsville

Multiscale Geometric Methods in Astronomical Data Analysis Institute for Pure and Applied Mathematics, UCLA November 8-12, 2004

From Data to Astronomical Goals



Astronomical Goals

Exploratory analysis: little prior information – few assumptions

Use simplest possible nonparametric model: *piecewise constant* (allows exact calculation of marginalized likelihoods)

Take account of:

observational noise exposure variations arbitrary sampling, gaps point spread functions (1D, 2D+ in progress)

intrinsic

-constr

Just say no to:

pretty pictures; smoothing; continuous representations bins or pixels (unless raw data in this form) methods tuned for specific structures, *e.g.* beamlets resampling sensitivity to some global structures, *e.g.* periodic

Piecewise Constant Model (Partition)

Can simply ignore gaps – model says nothing about signal there.



Smoothing and Binning

Old views: the best (only) way to reduce noise is to smooth the data the best (only) way to deal with point data is to use bins

New philosophy: smoothing and binning should be avoided because they ... discard information

degrade resolution

introduce dependence on parameters:

- degree of smoothing
- bin size and location

<u>Wavelet Denoising</u> (Donoho, Johnstone, ...) multiscale; uses no explicit smoothing

Adaptive Kernel Smoothing

<u>Optimal Segmentation</u> (*e.g.* Bayesian Blocks) Omni-scale -- uses neither explicit smoothing nor pre-defined binning

The Problem: Signal/Density Estimation

<u>The data:</u>Data cells: points, event counts, measurements, *etc*.Noise: known distribution, not necessarily additive (*e.g.* Poisson)

Distributed over a data space:
1D: time series, sequential data, ...
2D: images, photon maps, star/galaxy catalogs, ...
3D: galaxy redshift surveys, energy-resolved photon maps xD: time & energy resolved photon maps; data mining

<u>Goals:</u> Nonparametric representation of the underlying signal or density

Easy analysis of signal structure (*e.g.* clusters in point density) Suppression of noise (using prior knowledge of noise statistics) Objective, automatic analysis of large data sets (data mining)

Signal Model: Piecewise Constant

Represent signal as constant over elements of a partition of the data space.

Optimize model by maximizing model fitness over all possible partitions.

_ Nonparametric

Few prior assumptions about the signal:
 prior on signal amplitudes
 prior on number of partition elements

_ No limitation on the resolution in the independent variable.

_ Representation, while discontinuous, is convenient for further analysis

Local structure, not global

(a) -1 -2 -3 -4 -5 -6 -7 -8 -9 -10 -11 -12 -13 -14 -15 -16 -17 -18 -19 -20 -21 -22 -23 -24 -25 -26 -27 -28 -29 -30 -31 -32 -



DATA CELLS: Definition

data space: the set of possible measurements in some experiment

data cell: a data structure representing an individual measurement

For a segmented model, the cells must contain all information needed to compute the model *cost function*.

In a specific application, the data cells may or may not:

- _ be in one-to-one correspondence to the measurements
- _ partition the entire data space
- _ overlap each other
- _ leave gaps between cells
 - contain information on adjacency to other cells

DATA CELLS: Event (Point) Data

Measurements:	Point coordinates
Data Space:	Space of any dimension
Signal:	Point density (deterministic or probabilistic)
Data Cell:	Voronoi cells for the data points

Suf. Statistics N = number of points in block V = volume of block

Max Likelihood: $(N / V)^N e^{-N}$ Posterior:N! (V - N)! / (V+1)!

Example: any problem usually approached with histograms (1D) positions of objects from a sky survey (2D) positions of objects in a redshift survey (3D)

DATA CELLS: Serial Measurements

Measurements: Values and error distribution of dependent variable at given values of independent variable, e.g. $X(t) \sim N(x,\sigma)$

Data Space: Interval, area, volume, ...

Signal: Variation of dependent variable

Data Cell: Measurement point

Suf. Statistics $x_n = x_n(t_n), t_n$, parameter(s) of error distribution: $\sigma_n, ...$

log posterior:

Example: time series, spectra, images, SOM output ...

DATA CELLS: Distributed Measurements

Measurement: Dependent variable

averaged over a range of independent variable

M

10

100

1000

 10^{4}

Data Space: Space of any dimension

Signal: Physical variable

Data Cell: Measurement and its interval

Suf. Statistics: x, σ_x , W(t) = window function

Posterior: see Bretthorst, G-S orthogonalization

Example: spatial power spectra of CMB

Blocks

Block: a set of data cells

Two cases: _ Connected (can't break into distinct parts) _ Not constrained to be connected

Model = set of blocks

Fitness function:

F(Model) = sum over blocks F(Block)

The Optimizer

```
best = []; last = [];
for R = 1:num_cells
[ best(R), last(R) ] = max( [0 best] + ...
reverse( log_post( cumsum( data_cells(R:-1:1, :) ), prior, type ) ) );
if first > 0 & last(R) > first % Option: trigger on first significant block
```

```
changepoints = last(R); return
```

end

end

```
% Now locate all the changepoints
index = last( num_cells );
changepoints = [];
while index > 1
changepoints = [ index changepoints ];
index = last( index - 1 );
end
```























Distribution of Galaxies in 3-space (and time)

Nature of the astrophysical process:

initial density fluctuations = (Gaussian?) random field
fluctuations grow over time due to gravity
end product = discrete galaxies

What is the best mathematical model of this process?



Distribution of Galaxies in 3-space (and time)

Structures observed in redshift surveys:

_ Voids (3D underdense cells) -- "voids"

Sheets (2D density excesses) -- "Zeldovich pancakes"

_ Nodes (1D density excesses) -- "classical clusters"

Optimum Partitions in Higher Dimensions

_ Blocks are collections of Voronoi cells (1D,2D,...) _ Relax condition that blocks be connected _ Cell location now irrelevant _ Order cells by volume Theorem: Optimum partition consists of blocks that are connected in this ordering $_$ Now can use the 1D algorithm, O(N²) _ Postprocessing step identifies connected block fragments

Data: Background





Data: Background + Source

Data: Background + Source



Data: Background + Source



Data: Voronoi Tessellation



Block Decomposition





































2-Point Joint Distribution

