#### Spectral fitting methods



### Welcome!



Probabilities of discrete events



Probabilities of hypotheses



Scientific application







### Welcome!



Probabilities of discrete events



Probabilities of hypotheses



Scientific application

Poisson distribution (by Abraham de Moivre) Bayesian inference (by Pierre-Simon Laplace)

Horse kicks (by Ladislaus Bortkiewicz)

Stigler's law of eponymy (Robert K. Merton)

## Welcome!

- Practical information
  - organisers
  - wifi
  - local information
  - rough agenda & goal
  - dinner
- First content block

#### Local information



## Agenda

#### • Morning

- Measurement process & statistics involved
- Background treatment
- Local best fits
- Lunch Cantine
- Afternoon
  - Global fits & probability distributions
  - Model comparison
  - Coffee 15:30
  - Combining information
  - Beyond X-ray spectra
  - Discussion & Questions
- End: 17:30. Dinner 19:00
- Morning
  - Extended sources, calibration
  - Discussion & Questions
  - Poisson knowledge & practical pointers

Goal: what methods exist what are their benefits & limitations what to pay attention to

## Dinner today 19:00

- U6+Bus 59
  - Dietlindenstr.
  - Bus 59 Giesing
  - Richard-Strauss-str.
- U6+U4
  - Odeonsplatz
  - U4 Arabellapark
  - Richard-Strauss-str.



## Dinner today 19:00

- U6+Bus 59
  - Dietlindenstr.
  - Bus 59 Giesing
  - Richard-Strauss-st
- U6+U4
  - Odeonsplatz
  - U4 Arabellapark
  - Richard-Strauss-st



## First block

- Overview & Introduction
  - Measurement process
  - Background & source regions
  - Linear algebra approximation
  - Likelihood & statistics
- after: background
- after: fitting

# What Counts

## Single spectral bin

- Bernoulli coin flip
  - k=0 (p)
  - k=1 (1-p)
- Binomial
  - n tries, first k successful
  - $-P=p^{k}(1-p)^{(n-k)}$
- Poisson
  - n $\rightarrow$ inf but pn= $\lambda$

rule of thumb: if n>20 and p<0.05 n>100 and np<10

$$P(X=k)=inom{n}{k}(p)^k(1-p)^{n-k}$$





## Single spectral bin

- Poisson
  - k: integer  $P(k) = e^{-\lambda} \frac{\lambda^k}{k!}$
  - $-\lambda$ : real (mean&variance)
  - Asymmetric
  - Integer
  - Positive
- Scaling
- Addition
- Subtraction

shape changes

λ

(Poisson distribution) Variability!

(Skellam distribution)

Samples Electronics (shot noise) Photon counting (Poisson noise)



## Single spectral bin

0.35

0.30

0.05

0.00

5

10

15

♀<sup>0.25</sup> × 0.20  $\lambda = 1$ 

 $\lambda = 4$ 

 $\lambda = 10$ 

20

- Poisson
  - k: integer  $P(k) = e^{-\lambda} \frac{\lambda^k}{k!}$
  - $-\lambda$ : real (mean&variance)  $_{0.10}^{0.15}$
- Gaussian
  - Mean (µ) & variance ( $\sigma^2$ ) =  $\lambda$
  - Mean (µ) & variance ( $\sigma^2$ ) = k
  - real, can be negative













## Approximation quality

- Tails have different slopes
  - Gauss high-end more permissive
  - Poisson low-end more permissive
- Right way: Poisson
- Historically: Gauss faster to evaluate

#### "Statistics"

Poisson

#### - Likelihood $\mathcal{L}(k|\lambda) = e^{-\lambda} \lambda^k / k!$ -2\*log $\rightarrow$ -2log $\mathcal{L}(k|\lambda) = 2\lambda - 2k \log \lambda + C$

• Gaussian

- Likelihood 
$$\mathcal{L}(k|\mu,\sigma) = \exp[-((x-\mu)/\sigma)^2/2]/\sqrt{2\pi\sigma^2}$$
  
-2\*log  $\rightarrow$  -2log  $\mathcal{L}(x|\mu,\sigma) = ((x-\mu)/\sigma)^2 + C$ 

#### "Statistics"

Poisson

- Likelihood 
$$\mathcal{L}(k|\lambda) = e^{-\lambda}\lambda^k/k!$$
  
-2\*log  $\rightarrow$  -2log  $\mathcal{L}(k|\lambda) = 2\lambda - 2k\log\lambda + C$   
CStat, Cash

Gaussian

- Likelihood  $\mathcal{L}(k|\mu,\sigma) = \exp[-((x-\mu)/\sigma)^2/2]/\sqrt{2\pi\sigma^2}$ -2\*log  $\rightarrow$  -2log  $\mathcal{L}(x|\mu,\sigma) = ((x-\mu)/\sigma)^2 + C$ Chi<sup>2</sup>

Does not mean they follow a chi<sup>2</sup> distribution!

Cash (1979)

## Multiple bins



Poisson

 $\mathcal{L}(k_1, k_2 | \lambda_1, \lambda_2) = e^{-\lambda_1} \lambda_{11}^k e^{-\lambda_2} \lambda_{22}^k$ 

 $2\lambda_1 - 2k_1\log\lambda_1 + 2\lambda_2 - 2k_2\log\lambda_2 + C$ 

Gaussian

 $\mathcal{L}(x_1, x_2 | \mu_1, \sigma_1, \mu_2, \sigma_2) = \exp[-((x_1 - \mu_1) / \sigma_1)^2] \exp[-((x_2 - \mu_2) / \sigma_2)^2]$ 

$$((x_1 - \mu_1)/\sigma_1)^2 + ((x_2 - \mu_2)/\sigma_2)^2$$

### Multiple bins



 $\overrightarrow{\lambda} = \overrightarrow{F} \cdot \underline{R}$  $C = 2\overrightarrow{\lambda} \cdot \overrightarrow{\lambda} - 2\overrightarrow{k} \cdot \log \overrightarrow{\lambda}$ 

# Backgrounds

## Backgrounds





Assume time, location-independence

 $k_{_{Src}},\!\lambda_{_{Src}},\!t_{_{src}},\!A_{_{Src}}$  $k_{bkg}, \lambda_{bkg}, t_{bkg}, A_{bkg}$ 



$$\vec{\lambda}_{\rm src} = \vec{F}_{\rm src} \cdot \underline{R}_{\rm src} + \vec{F}_{\rm bkg} \cdot \underline{R}_{\rm src}$$
$$\vec{\lambda}_{\rm bkg} = \vec{F}_{\rm bkg} \cdot \underline{R}_{\rm bkg}$$
$$C = 2\vec{\lambda}_{\rm src} \cdot \vec{\lambda}_{\rm src} - 2\vec{k}_{\rm src} \cdot \log \vec{\lambda}_{\rm src}$$
$$2\vec{\lambda}_{\rm bkg} \cdot \vec{\lambda}_{\rm bkg} - 2\vec{k} \cdot \log \vec{\lambda}_{\rm bkg}$$

Assumptions:

- area energy-independent - rate constant with area, time, location

Remember:  $\lambda$ =number / cm<sup>2</sup> / s / keV \* dE \* dt \* dA



 $\lambda$ =number / cm<sup>2</sup> / s / keV \* dE \* dt \* dA

## Background + Source

- src+bkg Gauss  $\rightarrow$  Gauss (subtractable, flats/darks)
- src+bkg Poisson  $\rightarrow$  Poisson
  - High counts (>100) in every single src and bkg bin  $\rightarrow$  Gauss + Subtract with bkg variance propagation
  - Subtract & model with Skellam distribution
  - Do the right thing and model <u>both</u> as Poisson
    - •
    - •
    - •

## Background + Source

- src+bkg Gauss → Gauss (subtractable, flats/darks)
- src+bkg Poisson  $\rightarrow$  Poisson
  - High counts (>100) in every single src and bkg bin  $\rightarrow$  Gauss + Subtract with bkg variance propagation
  - Subtract & model with Skellam distribution
  - Do the right thing and model <u>both</u> as Poisson
    - Poisson estimate of rate in each bin, independently
    - Function approximation of background
      - In counts (empirical model)
      - Physical background flux model
      - Fit simultaneously with source
      - Fit background model first, use best-fit background shape for source fit



 $\lambda$ =number / cm<sup>2</sup> / s / keV \* dE \* dt \* dA

## eROSITA background



https://wiki.mpe.mpg.de/eRosita/ScienceRelatedStuff/Background

#### Semi-physical background models



 $\rightarrow$  especially important for extended source

## Empirical background models

Maximize poisson likelihood at all bins →shape

Pros:

- Can contain physical knowledge & smoothness
- Small uncertainties
- 0 bin counts ok

Cons:

- Need to specify modelFit can be poor



Chandra

(XMM, Chandra, Swift models in-house)

## Empirical background models



Automated shape finding Simmonds, Buchner et al. (2017)

XMM/PN,MOS, Chandra/ACIS, NuSTAR, Suzaku, RXTE, Swift/XRT





Estimate most likely background rate in each bin

Add scaled to source region counts

(wstat, Xspec default if set to cstat with no background model) pgstat Pros:

 no model specification needed

#### Cons:

- no continuity
- unnecessarily large uncertainties
- need >0 counts per bin
# Inference with likelihoods

$$\mathcal{L}(\overrightarrow{k}|\theta_1, \theta_2, \dots, \theta_d, M, R, B, \dots)$$

Higher L: model under these parameters often makes this data Lower L: less frequently

 $\rightarrow$  Frequency of data  $P(D|\theta)$ 

Likelihood function at D, at parameter values (not a density)

#### Inference desiderata

• Parameter ranges allowed or probable (L, T, ..., physical parameters)  $P(\theta|D)d\theta$ 

Probability density

In infinitely small region: zero probability



#### Conditional probabilities

- Bayes theorem
- P(A|B) != P(B|A)
- Normalisation
- Parameter inference
- Model inference
- Interpretation

#### Conditional probabilities

$$egin{aligned} A &\frown B &= B &\frown A \ p(A &\frown B) &= p(B &\frown A) \ p(A|B)p(B) &= p(B|A)p(A) \ p(A|B) &= rac{p(B|A)p(A)}{p(B)} \ p( heta|D) &= rac{p(D| heta)p( heta)}{p(D)} \end{aligned}$$
 Bayes theo



# Parameter space exploration

#### Parameter space exploration

- Local optimization
  - LM, simplex, ... (many)
  - Monte carlo optimization
- Local sampling: MCMC
  - Tempering
  - Limitations
- Global optimization
  - Genetic algorithms (DE)
- Global sampling
  - Nested sampling



#### Best fit parameters



If many data are created under  $\theta$  logL interval -1 below best fit Contains true value 68% of realisations

Confidence interval

What was the question again? Are conditions fulfilled? What do unequal "errors" mean?

- If away from boundary
- If model is linear
- If ndata  $\rightarrow$  high (symmetric, single gauss)
- If  $\theta$  is true parameter

→ then

#### Best fit parameters



#### Calibrate a Confidence interval

#### Detection



#### Best fit distributions



Convolution of

True parameter distribution + Measurement error & analysis method

Confidence intervals

Histogram of best fits

Meaning? Upper limits?

Cumulative distribution

Clean solution: Model population distribution (HBM) Buchner+17a

# Sampling





For example with a grid

### Posterior grid

- evaluate *likelihood* at every point
  - how prone is the process to produce the observed data
- Compute relative importance:



• Grab those that make up 90% of  $\sum \mathcal{L}$ 

- $Z=\mathcal{L}$  "evidence" is average likelihood

#### Posterior grid

- Result is dependent on placement
- Equal spacing in  $heta_1$  or in  $\log heta_1$ .
- Choice of spacing is called "prior"
- coin = investment in computing there, put coins where it is worthwhile



#### Bayesian posterior

parameter solutions weighted by their probability



Ρ

#### Credible intervals

Definitions:

Density  $\rightarrow$  cumulative  $\rightarrow$  quantiles

Highest Density Intervals

Borders (upper limits)



- Compare two parameter spaces by  $\sum \mathcal{L}\Big|_{M1} / \sum \mathcal{L}\Big|_{M2}$
- How many coins to put in M1, M2?
- model prior

#### Parameter Estimation vs. Model Comparison

- Remove coins contributing less than 10%.
- Under Bayesian inference, same problem:
  - comparing bags of hypotheses



- prior is measure, rule of averaging, deformation of space to "natural variables", investment in/weighting of sub-regions
- most common priors: uniform, log-uniform.
- model priors are relative size of spaces

## Curse of dimensionality

- kd grid  $\rightarrow$  infeasible
- Sample  $\theta$ 
  - $\boldsymbol{\theta}_1 \ \boldsymbol{\theta}_2 \ \boldsymbol{\theta}_3 \ \dots$
  - $W_1 W_2 W_3 ....$

(Posterior chains)

- Techniques:
  - Importance sampling
  - MCMC
  - Nested sampling

## Using posterior chains

- Posterior chain  $\theta_1 \theta_2 \theta_3 \dots$
- Find regions with high prob
- Compute prob. of regions

 $og_{10}L_{2-10keV}$ 

40

35

30 L

2

- Posterior predictions
- Derived quantities



#### Importance sampling

Ρ

Draw from proposal distribution Q

Weigh by  $Q(\theta)/P(\theta|D)$ 

 $\rightarrow$  weighted  $\theta$  chain

Advantages:

θ

- Efficient in low-d
- Parallelisable
- Can integrate parameter space

Disadvantages

- Need to find good proposal (VB)
  Poor scaling to 10-20d
  Poor performance if proposal is bad (variance indicator)

#### Markov Chain Monte Carlo L Starting point θ Х Loop forever: $\theta' = Normal(\theta, sigma_p)$ if $P(\theta'|D)/P(\theta|D) > U()$ : $\theta = \theta'$ add $\theta$ to chain θ

- Missing ingredient: transition kernel
- tune to the problems
- Fraction of visits ~ converges to ~ probability of hypothesis
- Where does chain spend 90% of its visits



### MCMC





# MCMC proposals

- Metropolis + Random Walk
- Goodman-Weare (emcee)
- HMC (Hamiltonian Monte Carlo)

 $\rightarrow$  animation

https://chi-feng.github.io/mcmc-demo/app.html

Random walk, HMC

# MCMC proposals

- Metropolis Random Walk
  - Adv: simple
  - Disadv: poor mixing
- Affine-invariant ensemble Goodman & Weare (2010)
  - Adv: auto-tuning for gaussian L
  - Disadv: poor mixing in bananas, collapses in high-d (Huijser+15)
- HMC (Hamiltonian Monte Carlo)
  - Adv: tunes itself to surface
  - Disadv: need gradients of models

# MCMC stopping

- MCMC theory: n→inf
- Trace plots
- Autocorrelation length
- Convergence tests
  - Detect if unreliable
  - Gelman-Rubin diagnostic
  - (many more)



Phases:

Identification Mixing

# Global optimization



#### Escaping local maxima: strategies

- Multiple random start positions
  - Augment local techniques
- Make surface easier
  - Tempering/Annealing
- Walker population
  - GW
  - Genetic algorithms (DE)



#### Genetic algorithms

**Mutation** 

Cross-over

Initial population







Selection



With fitness function (here: L)

New generation



#### Genetic algorithms

Differential evolution





#### Zooms into highest regions

moncar

Advantages:

- Global
- Robust to degeneracies, auto-tuning proposal
  Works in high-d & non-continuous parameters

Disadvantages:

- Some tuning parameters
  stopping criterion may not be meaningful
  Does not sample (only best-fit)

# Model comparison

## Model comparison

Buchner+14

- Empirical models
  - Information content
  - Prediction quality
- Component presence
  - Regions of practical equivalence
- Physical effects
  - Bayesian model comparison
  - Priors often well-justified



Betancourt (2015)

#### Information criteria

Akaike information criterion

```
Akaike (1973)
```

• Is more complex worth storing?

 $AIC = 2 * d - 2 * L_{max}$ AIC = 2 \* d + CStat

Advantages:

- rooted in information theory
- independent of prior

Disadvantages:

- No uncertainties, thresholds unclear

- ...



- Compare two parameter spaces by  $\sum \mathcal{L}\Big|_{M1} / \sum \mathcal{L}\Big|_{M2}$
- How many coins to put in M1, M2?
- model prior

## Punishing prediction diversity

(not number of parameters)



L high, V tiny

L medium, V medium

#### What to do with Z

• Z1, Z2

# $\frac{p(M1|D)}{p(M2|D)} = \frac{Z1 \cdot p(M1)}{Z2 \cdot p(M2)}$

Posterior odds ratio

Bayes factor od


#### • Z1, Z2



• model priors: leave to reader or motivated by theory

- Discard highly improbable model or marginalise
- Does  $rac{p(M1|D)}{p(M2|D)} = 3/1$  mean M2 is correct in a quarter of the cases?

# Global sampling

#### nested sampling idea

- MCMC: only consider likelihood ratios. Integration by vertical slices
- nested sampling: compute geometric size at various likelihood thresholds
- orthogonal, unique re-ordering of volume by likelihood





### nested sampling algorithm



- Start with volume 1, draw randomly uniformly 200 points
- remove one, volume shrinks by 1/200.





- draw a new one excluding the removed volume
- Unique ordering of space required: via likelihood

draw a new uniformly random point, with higher likelihood (the crux of nested sampling)

- Scanning up vertically, done at some point
- converges (flat at highest likelihood)

## Missing ingredients

- MCMC: Insert tuned transition kernel
- NS: Insert constrained drawing algorithm
  - General solutions: MultiNest, MCMC, HMCMC, Galilean, RadFriends, PolyChord

#### RadFriends / MultiNest

- Use existing points to guess contour
- Expand contour a little bit
- Draw uniformly from contour
- Reject points below likelihood threshold
- RadFriends: Compute distance at which every point has a neighbor. Bootstrap (Leave out) for safety.
- MultiNest clusters and uses ellipses



Animation:

https://johannesbuchner.github.io/mcmc-demo/app.html#RadFriends-NS,standard (via chi-feng.github.io)

#### • Z1, Z2



model priors: leave to reader or motivated by theory

- Discard highly improbable model or marginalise
- Does  $rac{p(M1|D)}{p(M2|D)} = 3/1$  mean M2 is correct in a quarter of the cases?

• Z1, Z2

## $\frac{p(M1|D)}{p(M2|D)} = \frac{Z1 \cdot p(M1)}{Z2 \cdot p(M2)}$

Posterior odds ratio

Bayes factor od



#### • Z1, Z2



• model priors: leave to reader or motivated by theory

- Discard highly improbable model or marginalise
- Does  $rac{p(M1|D)}{p(M2|D)} = 3/1$  mean M2 is correct in a quarter of the cases?

## Calibrating model decisions

- Model probabilities  $\rightarrow$  decisions
- False decision rate (false positives/negatives)
  - Monte Carlo simulations (parametric bootstrap)

Buchner+14

## Calibrating model decisions



**False negatives** Non-decisions

Buchner+14

Advantages:

- Get rid of parameter prior dependences
- Have frequentist properties of Bayesian method
  Completely Bayesian treatment + decisions

**Disadvantages:** 

- Can be computationally expensive

## Frequentist properties of Bayesian methods

- Make decisions
  - Is parameter greater than C?
  - Is this model "better" than the other?
- Parametric bootstrap
  - Monte Carlo simulation allow arbitrary complexity

## Model comparison



## Agenda

- Yesterday:
  - Basic statistics, problem setup
  - Parameter estimation methods, credible & confidence intervals
  - Model comparison methods
  - Visualisations
- Today:
  - Extended sources, calibration
  - Stacking information
  - Discussion & Questions
  - Practical pointers & Wrap-up
- Not covered:
  - Tools
  - Pile-up & variability

Spectroscopy + lower E + higher E + imaging + time Outside spectroscopy

# Spectra with few counts

## Spectra with few counts

- Are nothing special
- Poisson likelihood + good background handling
- 0 counts

• Think in terms of allowed regions

L, N<sub>H</sub> from X-ray spectrum



### L, N<sub>H</sub> from X-ray spectrum



## Intrinsic parameter distributions

## Example

- Measurement gave
  - x1 = 4 + -1x2 = 5 + -0.1
- Generate 100 samples for each

100xN matrix of x



• Evaluate model

100xN matrix of x



100xN matrix of F



• Sum probabilities



- Multiply probabilities
- Then try out other model parameters









## Behaviour

• Generate from 6+-2 with measurement errors



 $\sigma$ 

















 $\sigma$




























## Practical pointers

## Practical advice

- You can do this in any package!
- State what you are doing
- CStat (Poisson)
- Background with functions (check fit)
- Visualise, visualise, visualise
- Show posterior distributions & fits in data space
- Vary priors & assumptions
- Use nested sampling, MCMC with care
- Make simulations
- Ask for help

isis, sherpa, spex, xspec, 3ml, ...

## Contact points for questions

- Ask a colleague
- Astrostatistics Facebook group
- XSPEC Facebook group
- @MPE
  - J Michael Burgess– Johannes Buchner