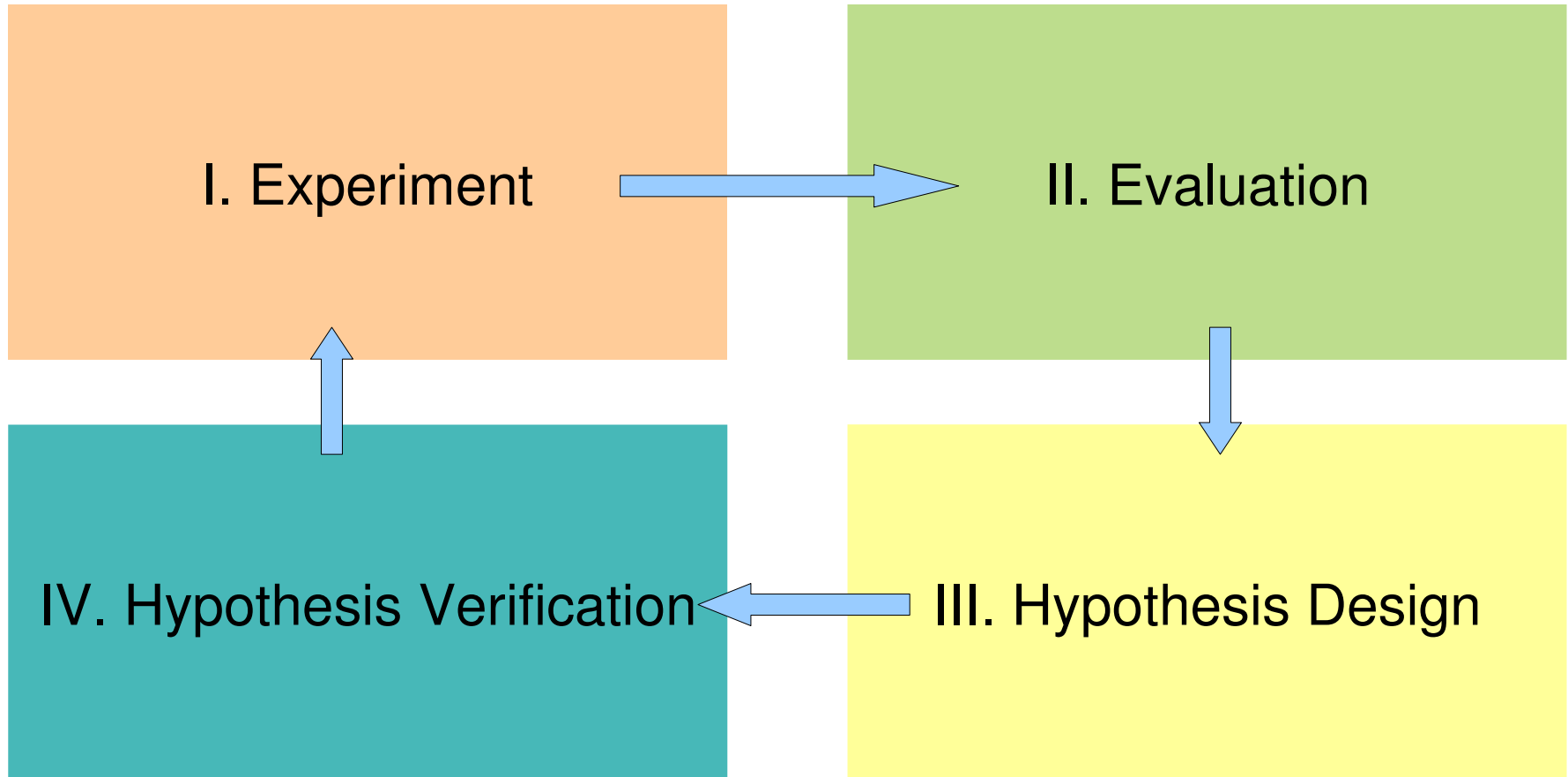


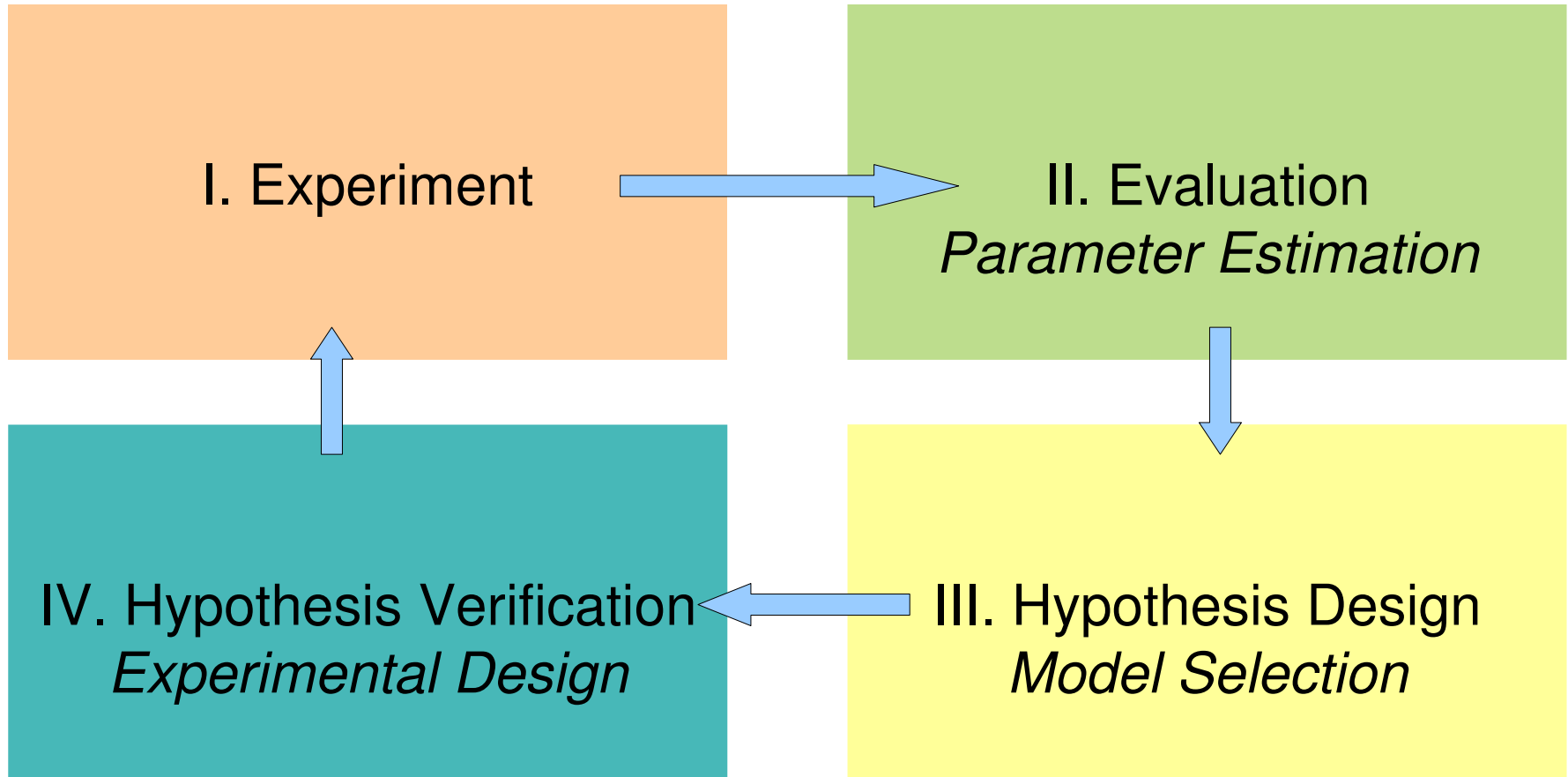
Bayesian Inference in Physics

U. von Toussaint

Scientific Inference Cycle



Scientific Inference Cycle



Prerequisite: Consistent reasoning...

I. Scientific Inference	Inference in Science Processing of Information
II. Model Comparison	Basic Concept Mass Spectroscopy
III. Experimental Design	Basic Concept Nuclear Reaction Analysis
IV. Numerical Interlude	Nested Sampling
V. Conclusion	Summary Outlook

Science: prior information + new data  new knowledge

Prior information:

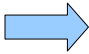
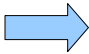
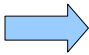
- old data
- calibration measurements, validation data
- theoretical considerations
- parameters
- model

continuous learning process 

New data:

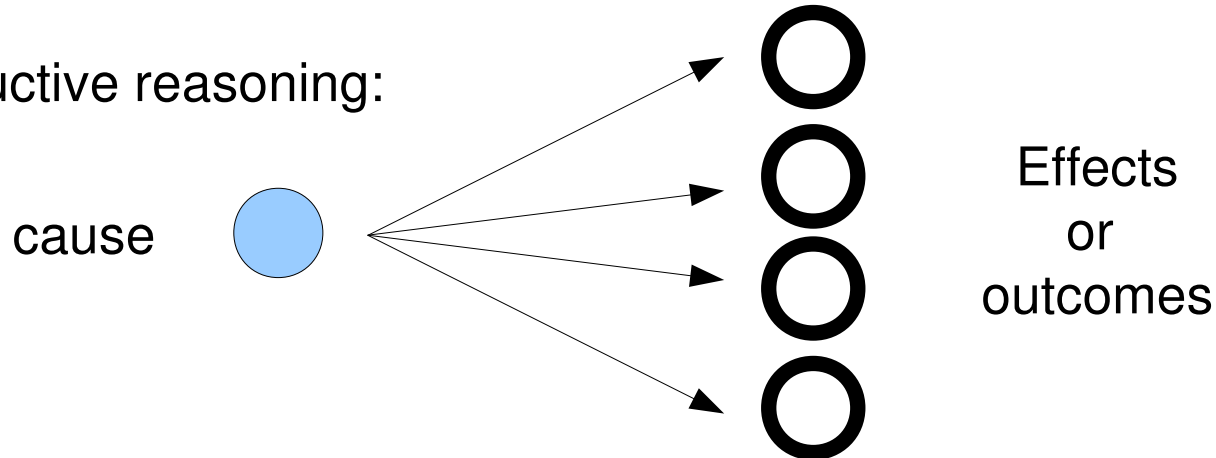
- measurements
- calibration
- theories

But: prior information and data are uncertain

-  new knowledge and derived hypotheses are uncertain
-  inductive (instead of deductive) reasoning required
-  **Calculus** for “uncertainties” needed (quantification)

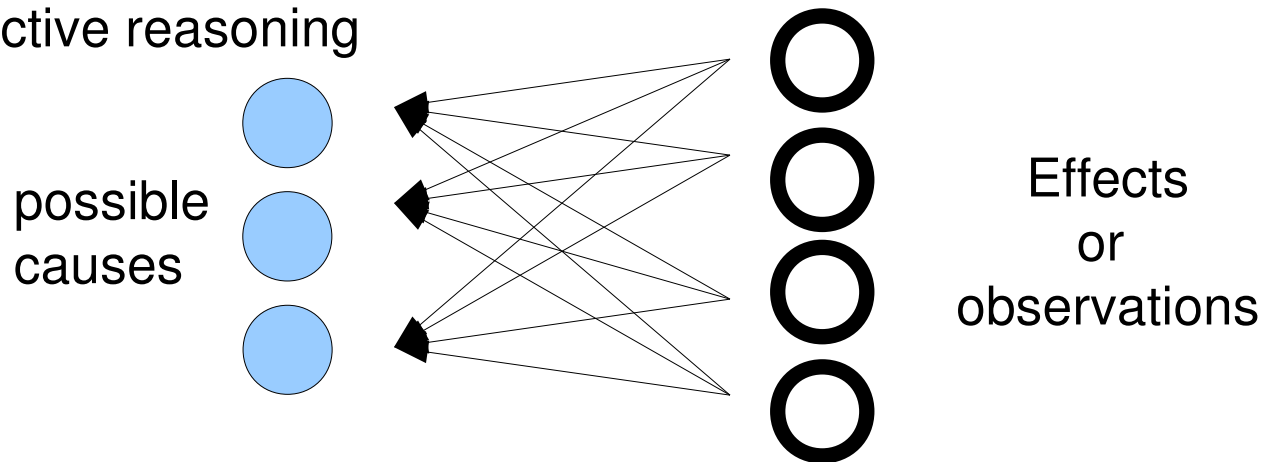
Science: prior information + new data  new knowledge

a) deductive reasoning:



Example: fair coin: $p(7 \text{ heads} \mid 10 \text{ tosses}) = ?$

b) inductive reasoning



Example: 7 heads out of 10 tosses: $p(\text{fair coin}) = ?$

Calculus:

Cox (1946): Basic requirements (i.e. transitivity, consistency) single out *usual rules of probability theory* to handle uncertainty

Please note: probability here *not restricted* to frequency interpretation:

Degree of belief about a proposition

Data:

- statistical (counting statistics)
- measurement uncertainty (ruler)
- systematic (e.g. misalignment)
- outliers

Hypotheses:

- parameter of interest
- nuisance parameters
- physical models
- future data

Notation: $p(x|I)$

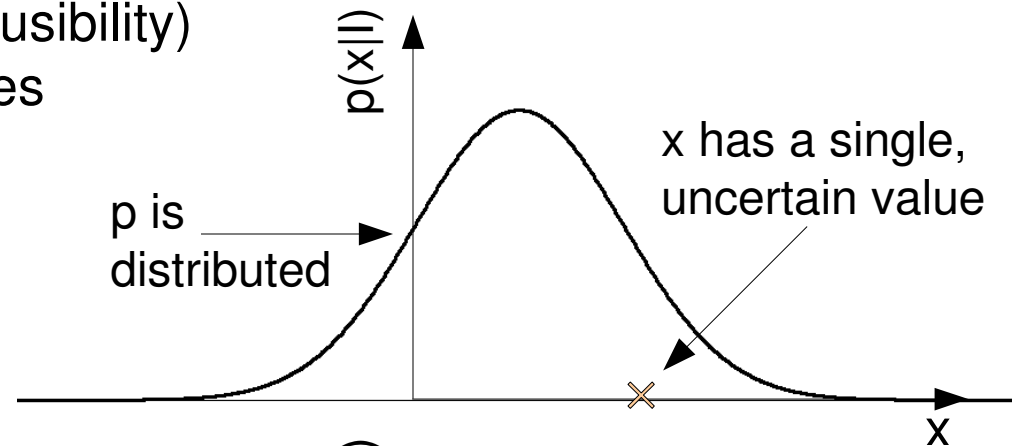
$p(x|I)$ describes how probability (plausibility) is distributed among possible choices for x for the case at hand (information I)

Calculus:

Cox (1946): Basic requirements (i.e. transitivity, consistency) single out *usual rules of probability theory* to handle uncertainty

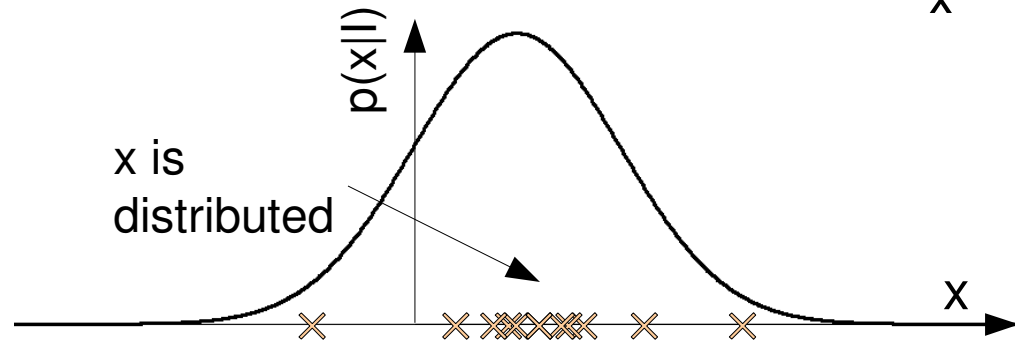
Bayesian interpretation:

$p(x|I)$ describes how probability (plausibility) is distributed among possible choices for x for the case at hand (information I)



Frequentist interpretation:

$p(x|I)$ describes how x is distributed throughout an infinite (hypothetical) ensemble
probability = frequency



Processing of information: Combination of (cond.) probability distributions

Conditional Probabilities:

$p(A|B)$: probability of proposition A given truth of proposition B:
quantification of uncertainty of A

Probability Theory Axioms:

- **sum rule** (OR): $p(H_1 + H_2 | I) = p(H_1 | I) + p(H_2 | I) - p(H_1, H_2 | I)$

- **product rule** (AND): $p(H_1, D | I) = p(D | H_1, I) p(H_1 | I)$
 $= p(H_1 | D, I) p(D | I)$



$$p(H_1 | D, I) = \frac{p(D | H_1, I) p(H_1 | I)}{p(D | I)}$$

**Bayes
Theorem**

Bayes' theorem:

$$p(H_1 | D, I) = \frac{p(D | H_1, I) p(H_1 | I)}{p(D | I)}$$

Posterior = $\frac{\text{Likelihood} \cdot \text{Prior}}{\text{Evidence}}$

Prior: knowledge before experiment (logically)

Likelihood: Probability for data if the hypothesis was true

Posterior: Probability that the hypothesis is true given the data

Evidence: normalization; important for model comparison



Maximum Likelihood approach (parameters which maximise probability for data) **does not** yield most likely parameters*!

Examples: $p(\text{wet street} | \text{rain}, I) \neq p(\text{rain} | \text{wet street}, I)$

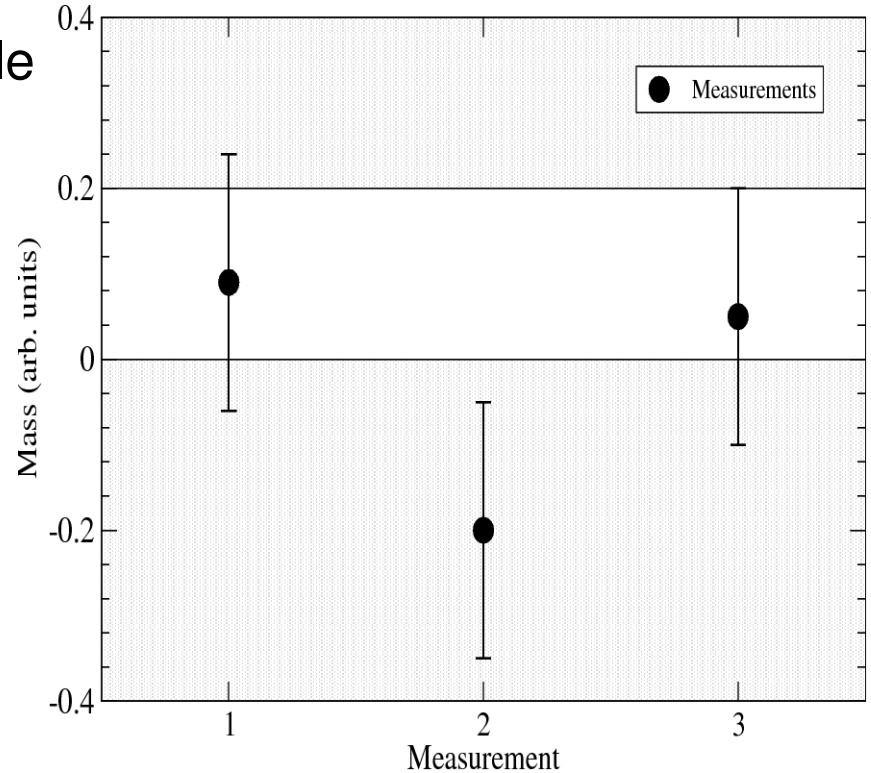
$p(\text{female} | \text{pregnant}, I) \neq p(\text{pregnant} | \text{female}, I)$

*) in general, except e.g. flat, unbounded priors

Processing of information

Toy Example: Mass of elementary particle

- Prior knowledge: $0 \leq m \leq m_{\text{upper}} = 0.2$
- measurement uncertainty $\sigma = 0.15$
- measured data:
 - $d_1 = 0.09$,
 - $d_2 = -0.2$,
 - $d_3 = 0.05$



Maximum likelihood: $m = -0.02$ (<0!)

Bayes: 1) Assign **prior**: $p(m | I) = \begin{cases} 1/m_{\text{upper}}, & \text{if } 0 \leq m \leq m_{\text{upper}}; \\ 0, & \text{otherwise,} \end{cases}$

2) Assign **likelihood**: $p(d | m, \sigma, I) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(d - m)^2}{\sigma^2}\right)$

3) Compute **posterior** using Bayes' theorem:

posterior :
$$p(m | \mathbf{d}, \sigma, I) = \frac{p(m | I) \prod_{i=1}^N p(d_i | m, \sigma, I)}{Z}$$

contains complete information:

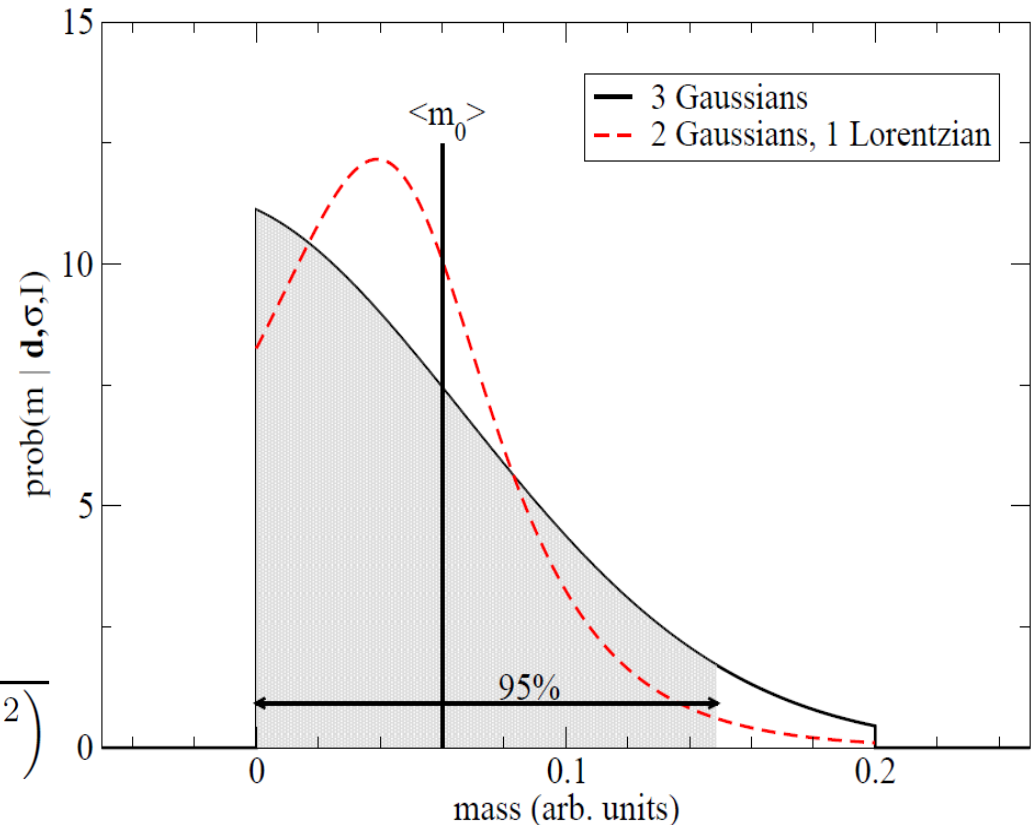
summarizing quantities:

- mode: $m=0$
- mean: $\langle m_0 \rangle = 0.06$
- 95%-interval: $[0; 0.145]$

Non-gaussian likelihoods: 😊

E.g.
$$p(d | m, \beta, I) = \frac{\beta}{\pi (\beta^2 + (d - m)^2)}$$

for d_3 with $\beta=0.05$



Reasoning about parameter a :

- (uncertain) prior information
- + (uncertain) measured data
- + physical model

$$p(a|I)$$

$$d=D \pm \varepsilon$$

$$D=f(a)$$

prior distribution

$p(d|a)$ likelihood distribution

+Bayes theorem

$$p(a | d) = \frac{p(d | a) \times p(a)}{p(d)}$$

posterior distribution

+an additional (nuisance) parameter b :

$$\begin{aligned} p(a | d) &= \int db p(a, b | d) \\ &= \int db \frac{p(d | a, b) p(a, b)}{p(d)} \end{aligned}$$

**Marginalization:
generalized error
propagation**

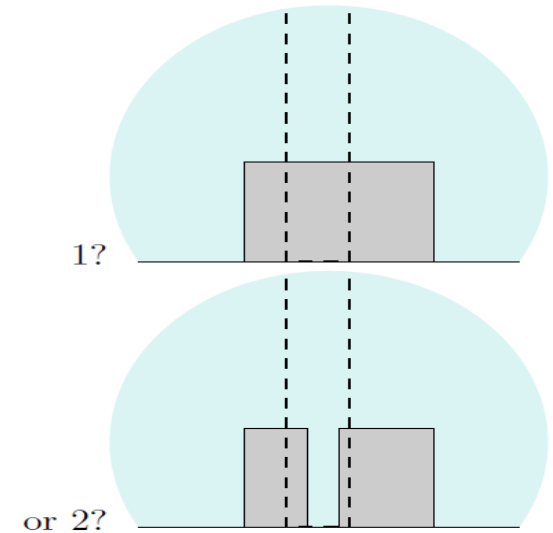
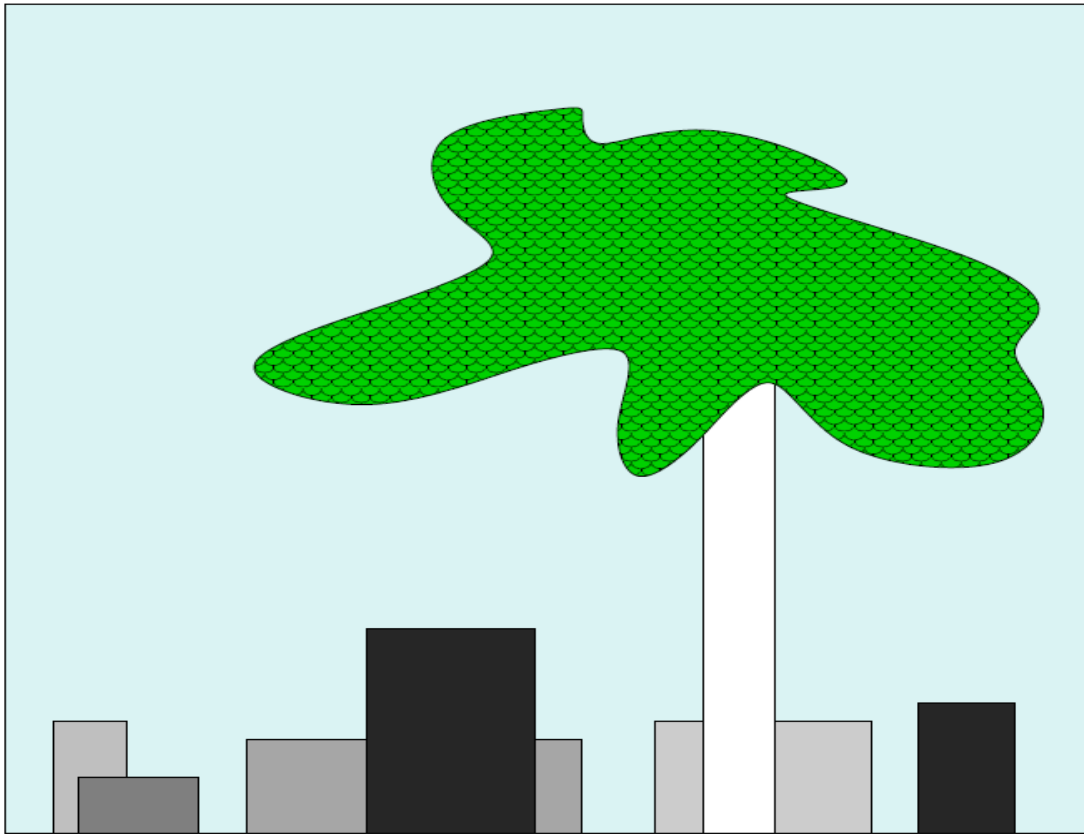
Clear recipe how to tackle a problem – possibly demanding mathematics/numerics

- 1) Quantify information at hand in probability distributions
- 2) Multiply probability distributions
- 3) Marginalize nuisance parameters
- 4) Analyze posterior distribution

II. Model Comparison

Basic Concept
Mass Spectroscopy

How many boxes are in the picture^{*)}?



Desired: **Occam's Razor** (Prefer simpler models (that fit the data))

^{*)}MacKay, Information theory

Model Comparison

Bayesian Approach:

$I = (M_1 + M_2 + \dots)$ — Specify a set of models.

$H_i = M_i$ — Hypothesis chooses a model.

Posterior probability for a model:

$$p(M_i|D, I) = p(M_i|I) \frac{p(D|M_i, I)}{p(D|I)}$$

$$\propto p(M_i) \mathcal{L}(M_i)$$

Posterior
for θ



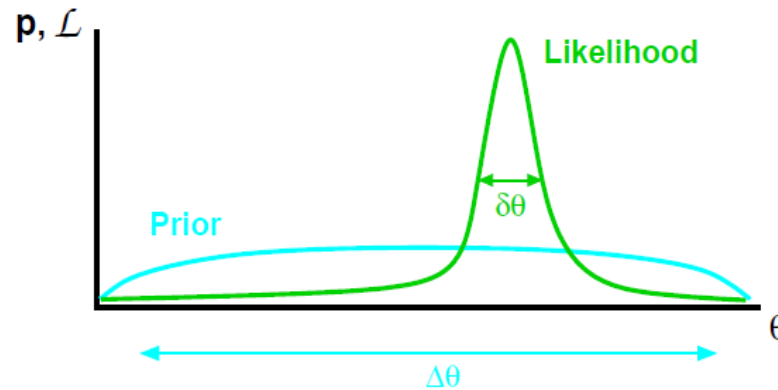
But $\mathcal{L}(M_i) = p(D|M_i) = \int d\theta_i \underbrace{p(\theta_i|M_i)p(D|\theta_i, M_i)}$.

Likelihood for model = Average likelihood for its parameters

$$\mathcal{L}(M_i) = \langle \mathcal{L}(\theta_i) \rangle$$

- Bayes Model Comparison requires always alternative models

The Occam Factor:^{*}



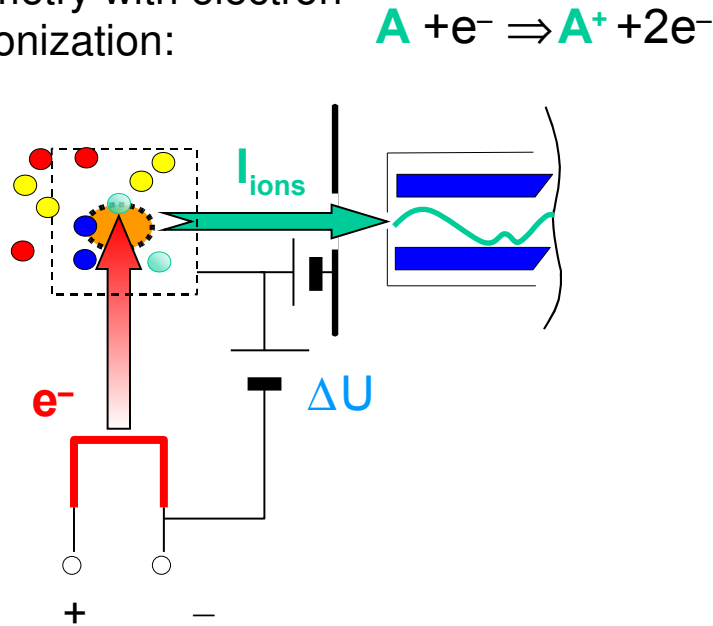
$$\begin{aligned} p(D|M_i) &= \int d\theta_i p(\theta_i|M) \mathcal{L}(\theta_i) \approx p(\hat{\theta}_i|M) \mathcal{L}(\hat{\theta}_i) \delta\theta_i \\ &\approx \mathcal{L}(\hat{\theta}_i) \frac{\delta\theta_i}{\Delta\theta_i} \\ &= \text{Maximum Likelihood} \times \text{Occam Factor} \end{aligned}$$

Models with more parameters often make the data more probable— *for the best fit*.

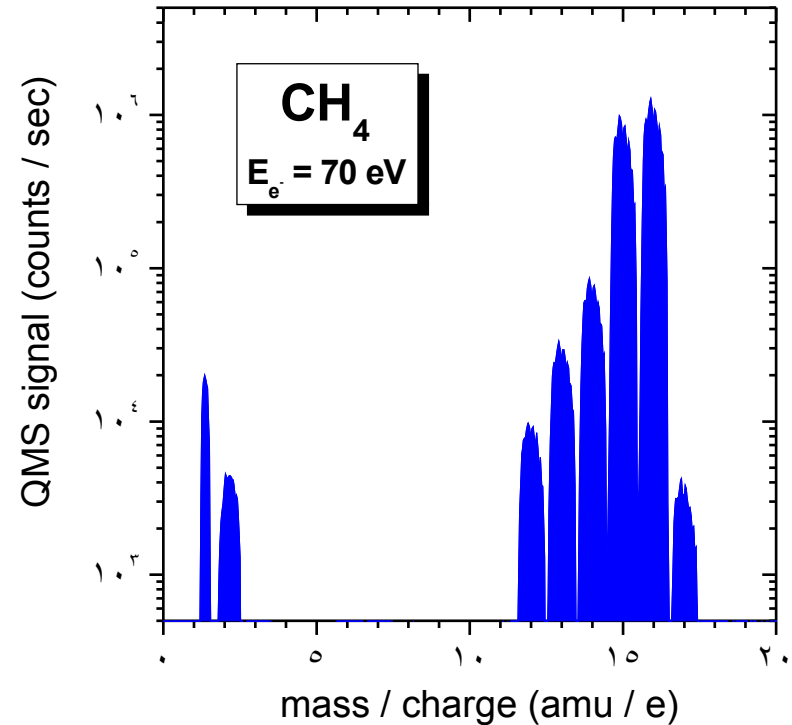
Occam factor penalizes models for “wasted” volume of parameter space.

^{*}) T. Loredo, Garching, 2003

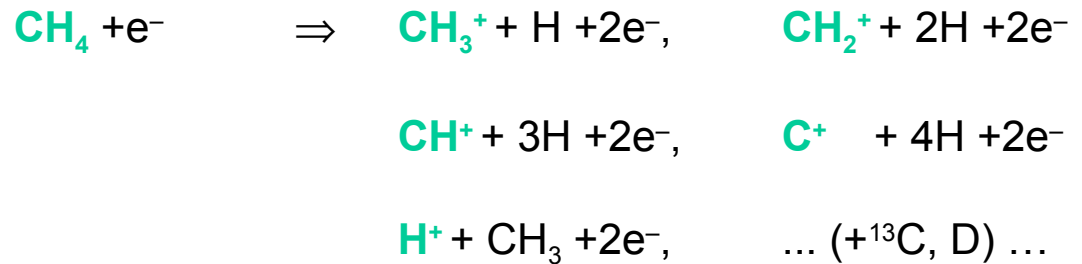
Quadrupole mass spectrometry with electron impact ionization:



Electron impact ionization can lead to complex, **instrument** dependent fragmentation patterns:



Versatile tool for neutral gas analysis: fast, high dynamic range, flexible,...





***"a mass spectrometrists is someone,
who figures out what something is,
by smashing it with a hammer
and looking at the pieces"***
Quote from Th. Schwarz-Selinger

Model: $\mathbf{d}_j = \mathbf{C}\mathbf{x}_j + \epsilon_j$ with cracking matrix \mathbf{C} (\mathbf{C} , \mathbf{x} are uncertain/unknown)

Likelihood: $p(\mathbf{D}|\mathbf{C}, \mathbf{X}, \{\mathbf{S}\}, E, I) =$
$$\prod_j \frac{1}{\prod_i \sqrt{2\pi s_{ij}}} \exp\left(-\frac{1}{2} (\mathbf{d}_j - \mathbf{C}\mathbf{x}_j)^T \mathbf{S}_j^{-1} (\mathbf{d}_j - \mathbf{C}\mathbf{x}_j)\right)$$

\mathbf{S} : measurement uncertainties

E : number of species

Prior terms:

Concentrations \mathbf{x} : depending on experiment (e.g. gas mixture)

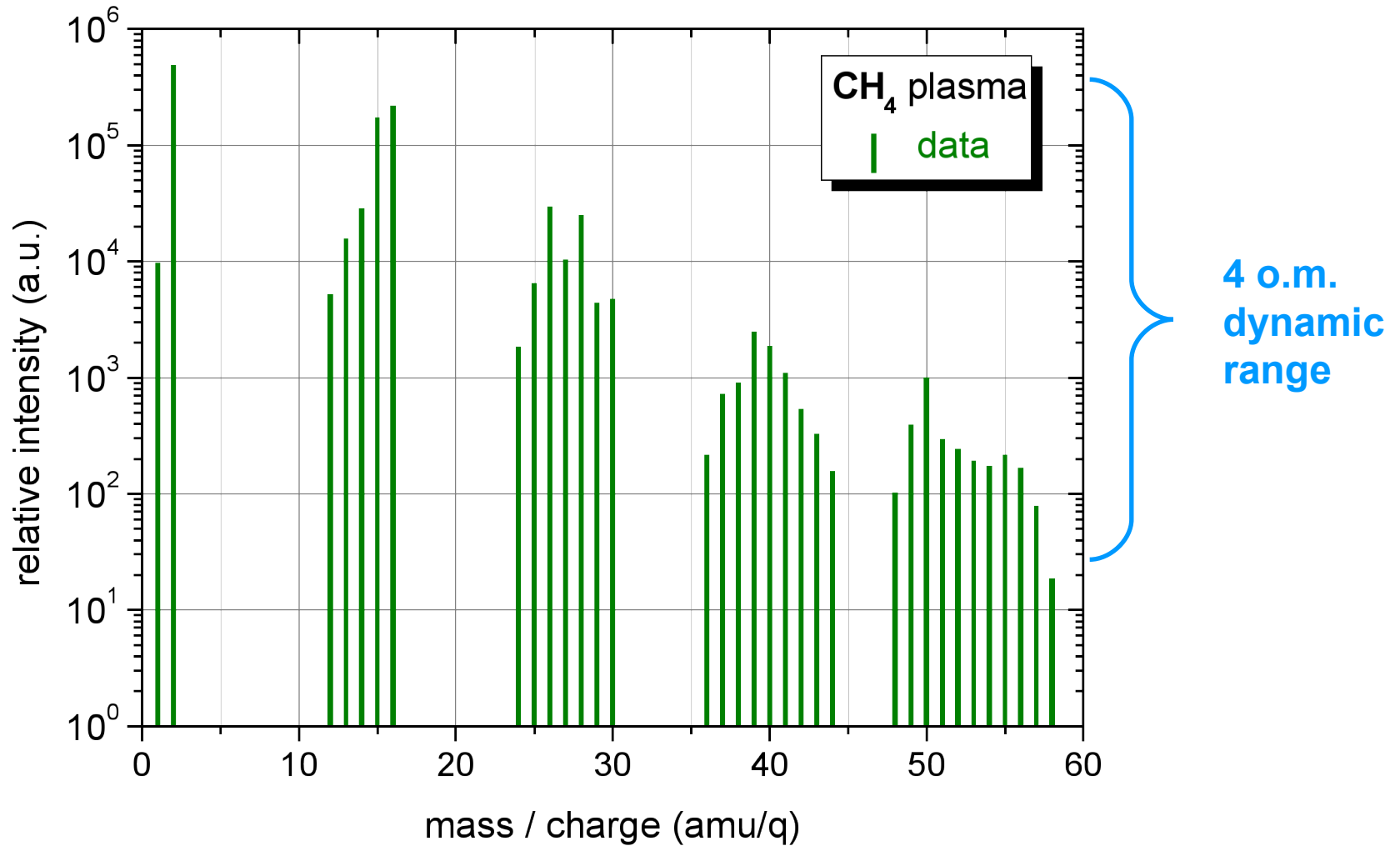
Cracking matrix \mathbf{C} : exponential prior based on point estimates of Cornu&Massot (1979)

Probability for a set of species $\{E\}$: $p(E|\mathbf{D}, \{\mathbf{S}\}, I) = \frac{p(E|I) p(\mathbf{D}|\{\mathbf{S}\}, E, I)}{p(\mathbf{D}|\{\mathbf{S}\}, I)}$

with $p(\mathbf{D}|\{\mathbf{S}, E\}, I) =$

$$\int d\mathbf{C} d\mathbf{X} p(\mathbf{C}|E, I) p(\mathbf{X}|E, I) p(\mathbf{D}|\mathbf{C}, \mathbf{X}, \{\mathbf{S}\}, E, I) \quad \leftarrow \text{High-dimensional integrals! MCMC-Integration}$$

Measured data:



input:

data:

signal + error of 34 mass channels
for 27 different plasmas conditions

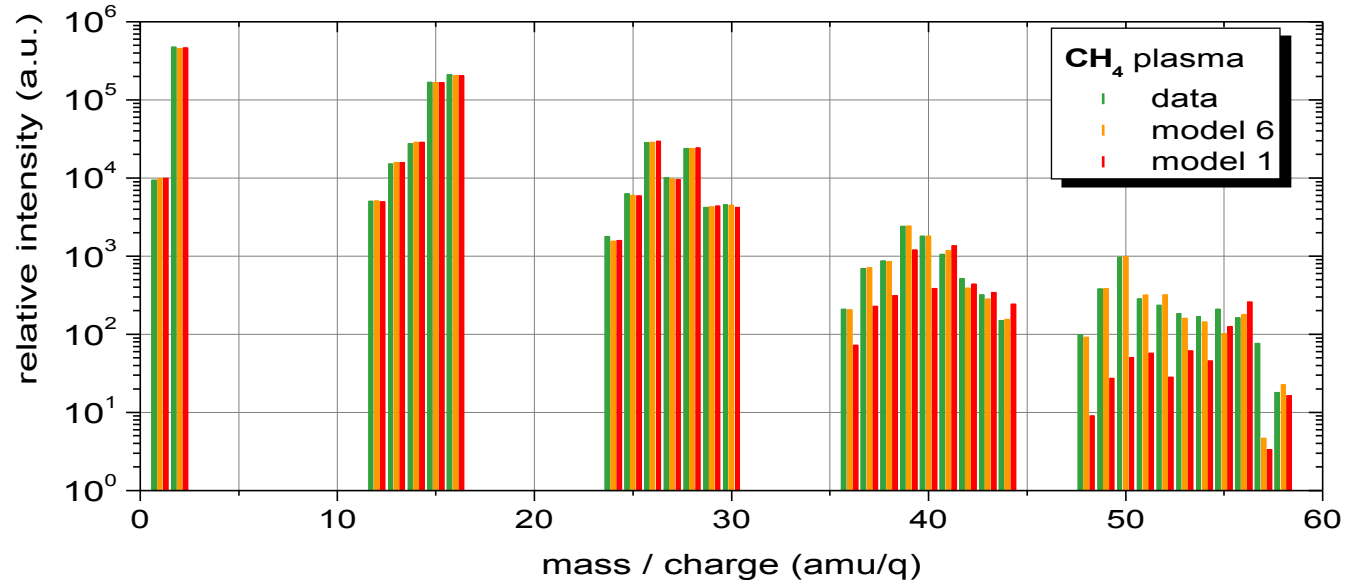
calibration measurements + error
for 11 species

prior:

cracking estimates for
14 species

output:

cracking pattern and
concentrations (+ errors!)
for 14 species
(e.g. C_4H_2 , C_3H_6)

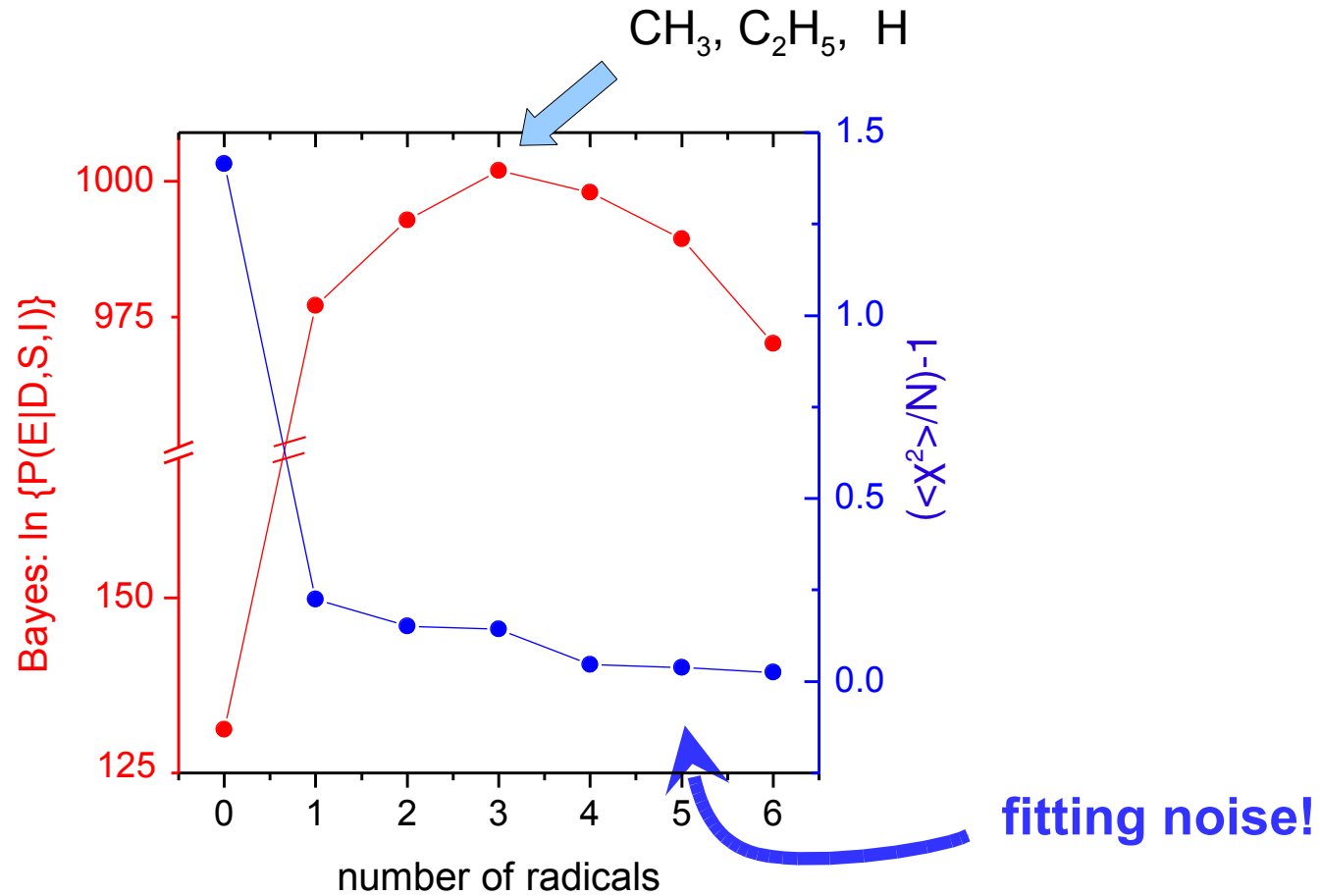


→ 'shifted' cracking pattern for radical CH_3 very bad

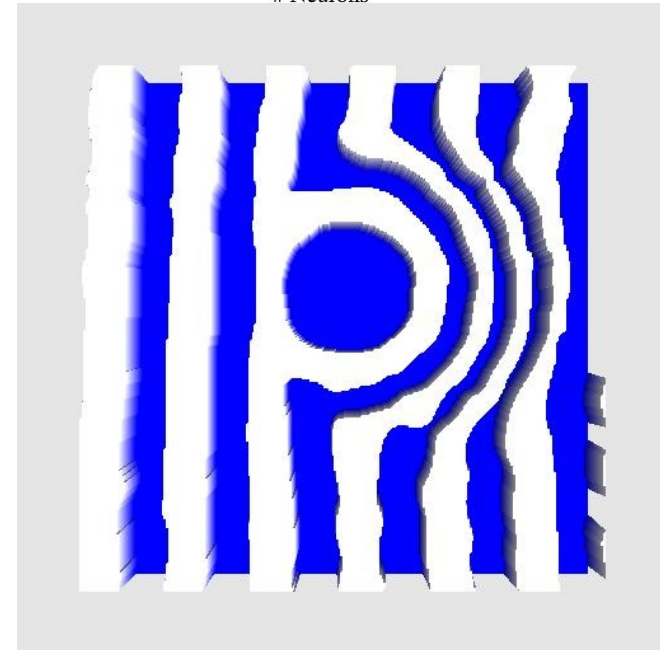
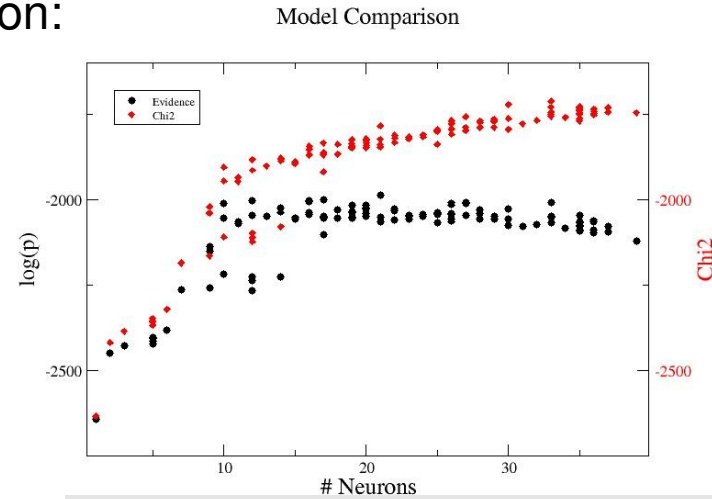
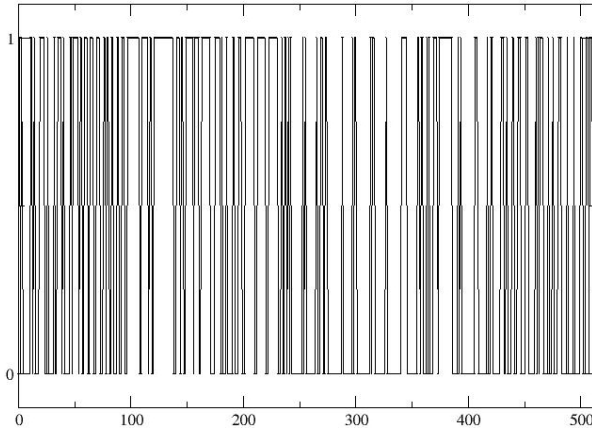
*) Th. Schwarz-Selinger, H. Kang, U. von Toussaint et al, various publications (2001-2007)

How many radicals are in the plasma?

model comparison with „**Occams Razor**“



An extreme example of model comparison:
Neural Networks without training data



U. von Toussaint et al, JAO, 2006

U. von Toussaint, MPE, 25.01.2012

III. Experimental Design

Basic Concepts
Nuclear Reaction Analysis

Design of experiments:

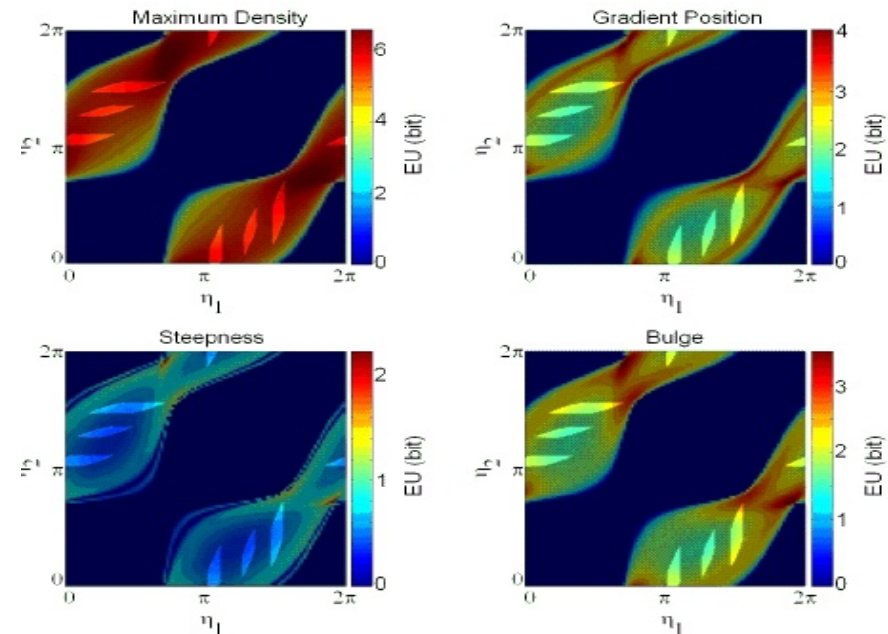
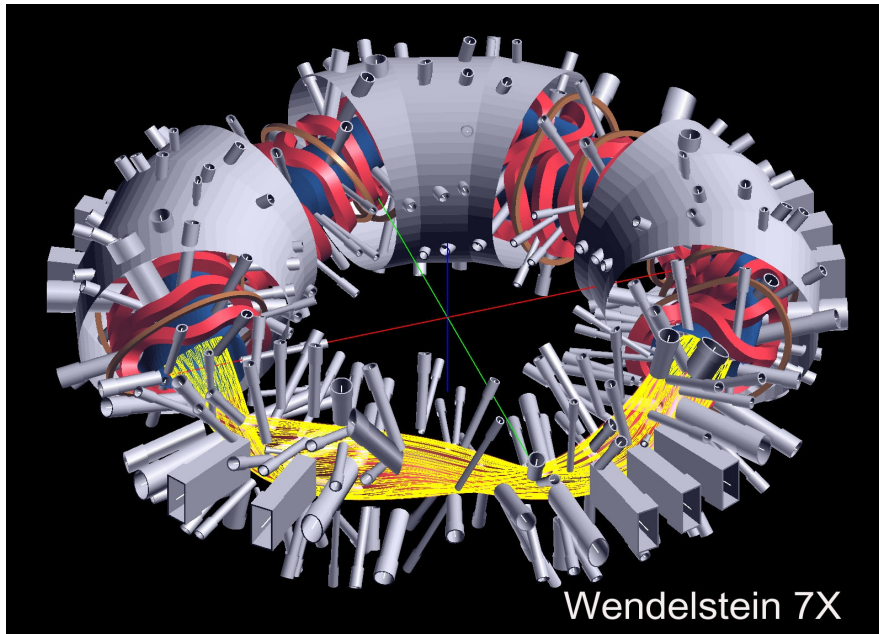
Maximize information gain of planned experiments

Best performance for various physical scenarios,
hardware constraints, parameter constraints, financial budget

How to quantify the strength (weakness) of an experiment?

How to quantitatively chose experimental design parameter(s) ?

(spectral bands, number and position of line-of-sights, measurement times)



Bayesian Decision Theory

Decisions depend on consequences

Might bet on improbable outcome if payoff is large

Utility functions

Compare consequences via utility quantifying the benefits of a decision



Choice of action: c (eg next accelerator energy E_0)

Utility= $U(c,o)$



Outcome: o (eg next yield $d(E_0)$)

Deciding amidst uncertainty

We are uncertain of what the outcome will be: average possible results

Expected Utility:
$$EU(c) = \sum_{\{outcomes\}} P(o|c) U(c, o)$$

Best choice maximizes EU: $c^* = \arg \max_c EU(c)$

Experimental Design

Information as Utility:

Goal: minimize estimation uncertainty for parameter(s) \mathbf{a} , maximize information gain

Utility function: Kullback-Leibler divergence K

(generalizes covariance-matrix)

$$Utility(d, c) = K(d, c) = \int da p(a|d, c) \ln \frac{p(a|d, c)}{p(a)} \quad (c=\text{action}, d=\text{expected data})$$

Putting it all together:

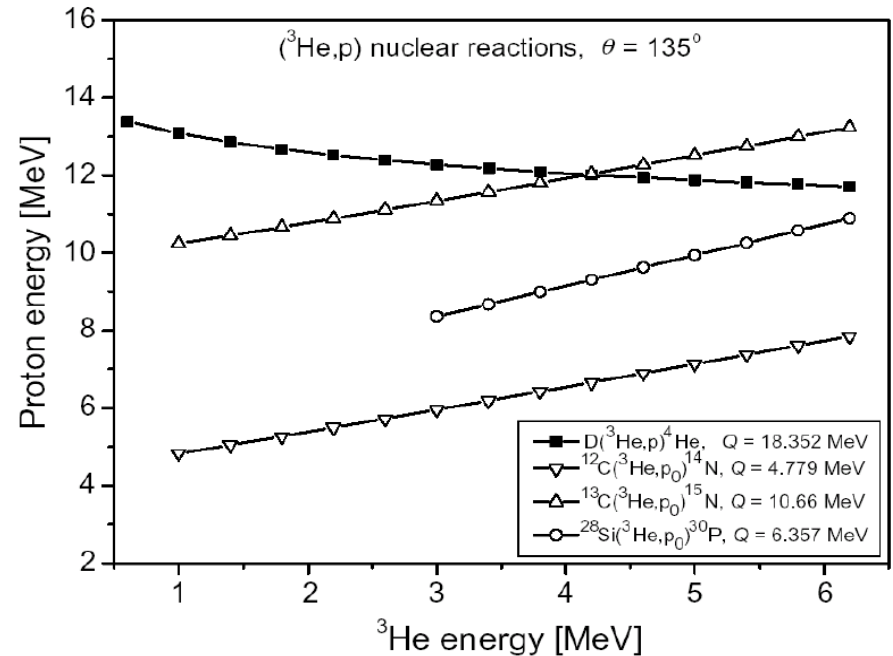
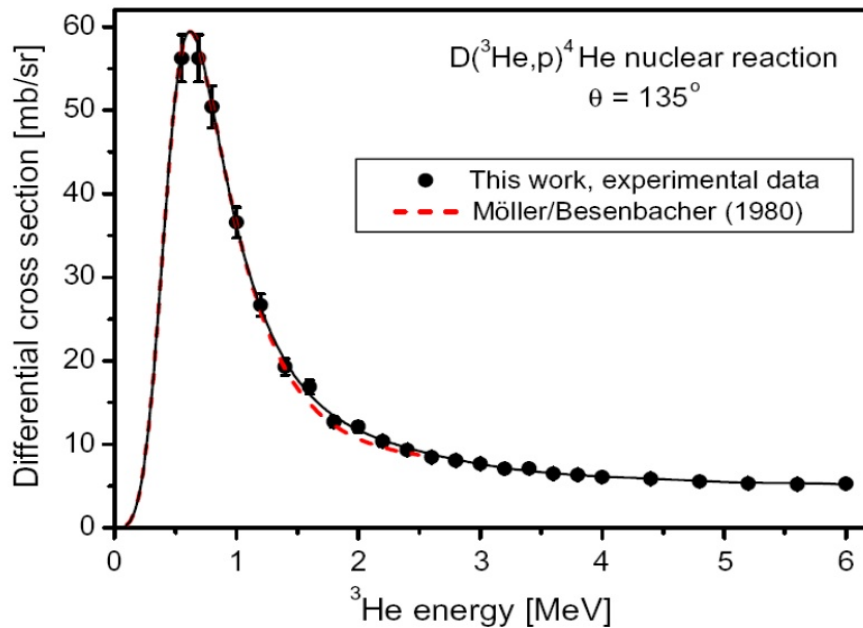
$$EU(c) = \int dd p(d|c) K(d, c) = \int da p(a) \int dd p(d|a, c) K(d, c)$$

Best choice maximizes EU: $c^* = \arg \max EU(c)$

High- dimensional integration necessary: Markov Chain Monte Carlo-Methods

- Hydrogen isotope depth profiling using NRA (not too many other ways)
- nuclear reaction: $d(^3\text{He},p)^4\text{He}$: background free, range

$$Y_i \propto \int_0^\infty dx \sigma (E (\vec{c}(x), x, E_{i0})) \left(\frac{dE}{dx} \right) \cdot c_j (x) + \epsilon_i$$



V. Alimov, J. Roth, NIMB, 2006

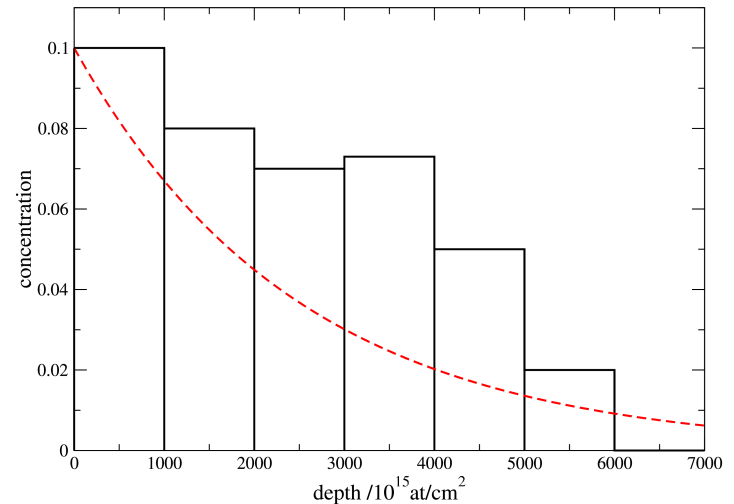
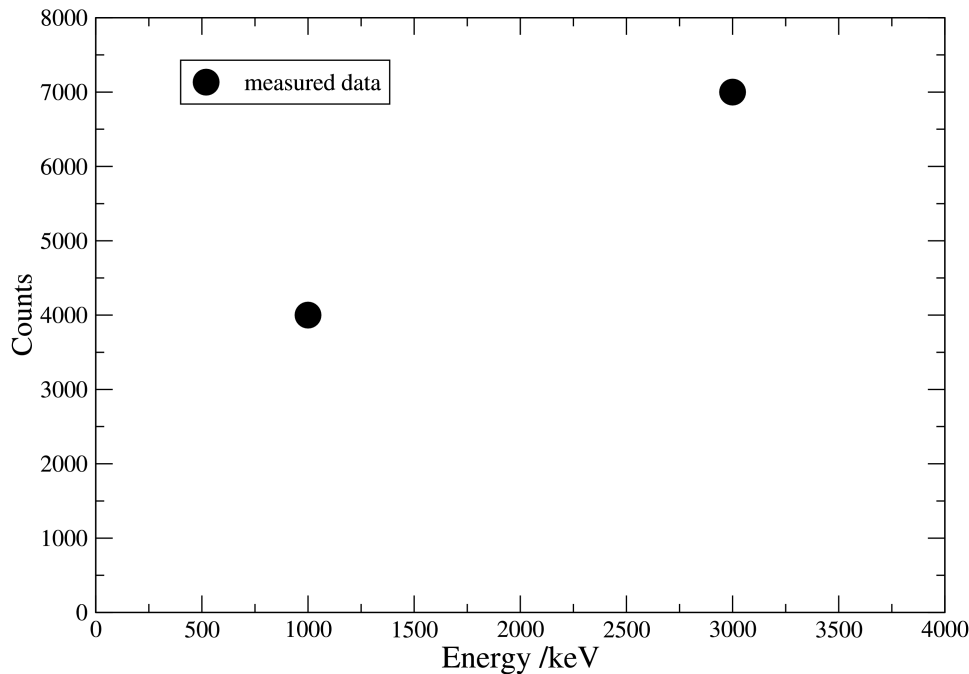
$$Y_i \propto \int_0^\infty dx \sigma (E (\vec{c}(x), x, E_{i0})) \left(\frac{dE}{dx} \right) \cdot c_j(x) + \epsilon_i$$

- depth profile has to be parametrized: piecewise constant or

analytically:

$$c(x) = \{a_0, a_1, a_2, a_3, \dots\} \quad \text{or}$$

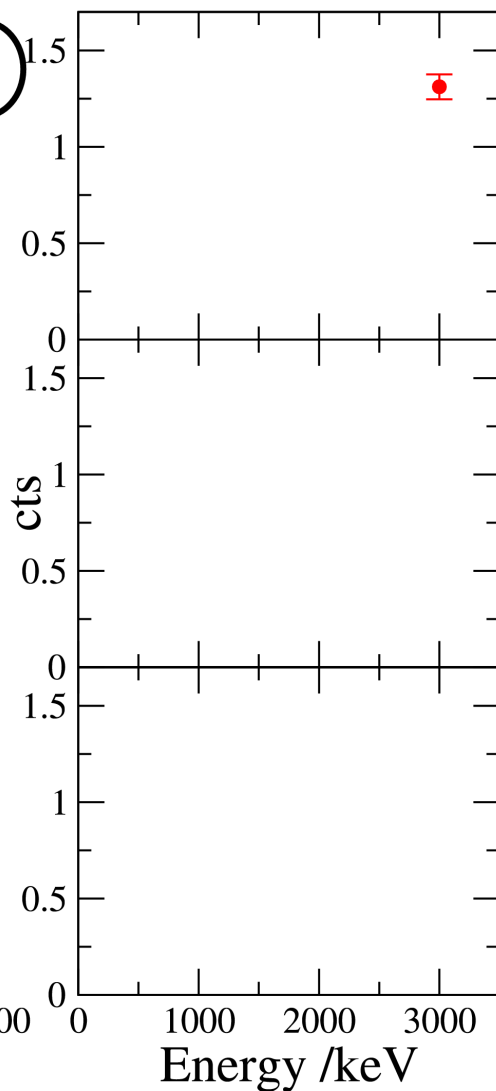
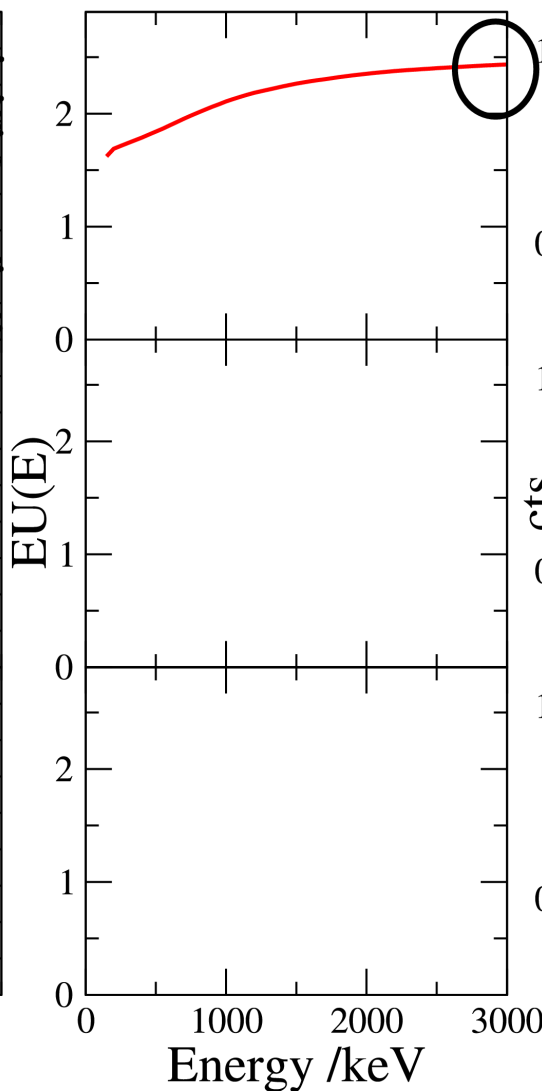
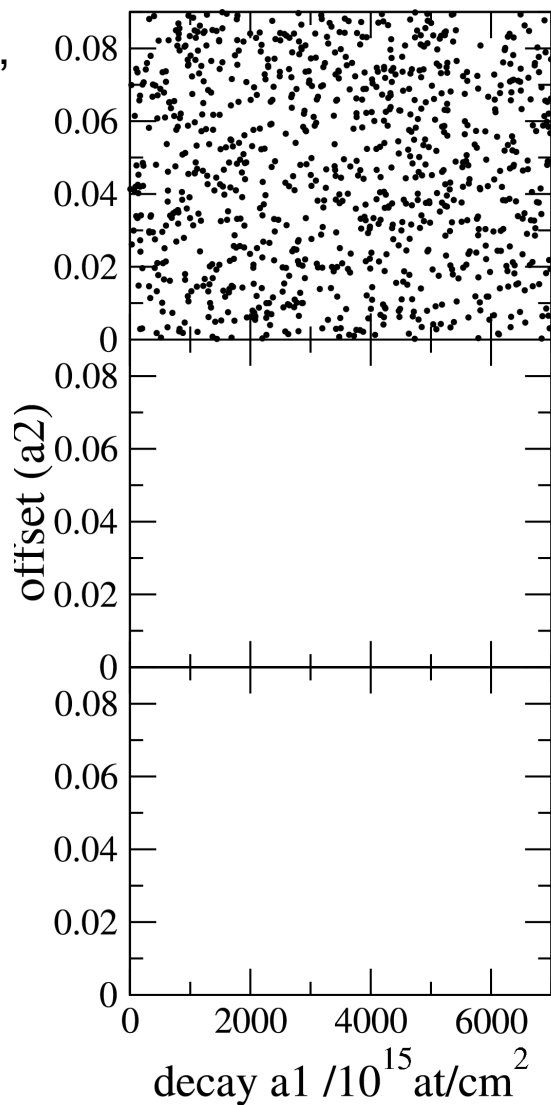
$$c(x) = a_1 \exp\left(-\frac{x}{a_2}\right) + a_0$$



Sequential Design: Which energy E_0 next?

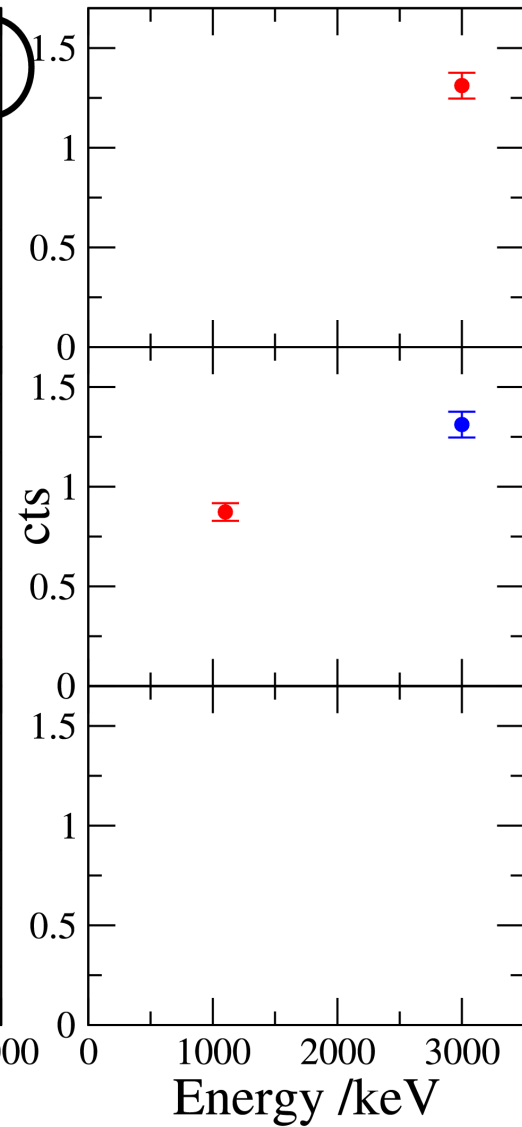
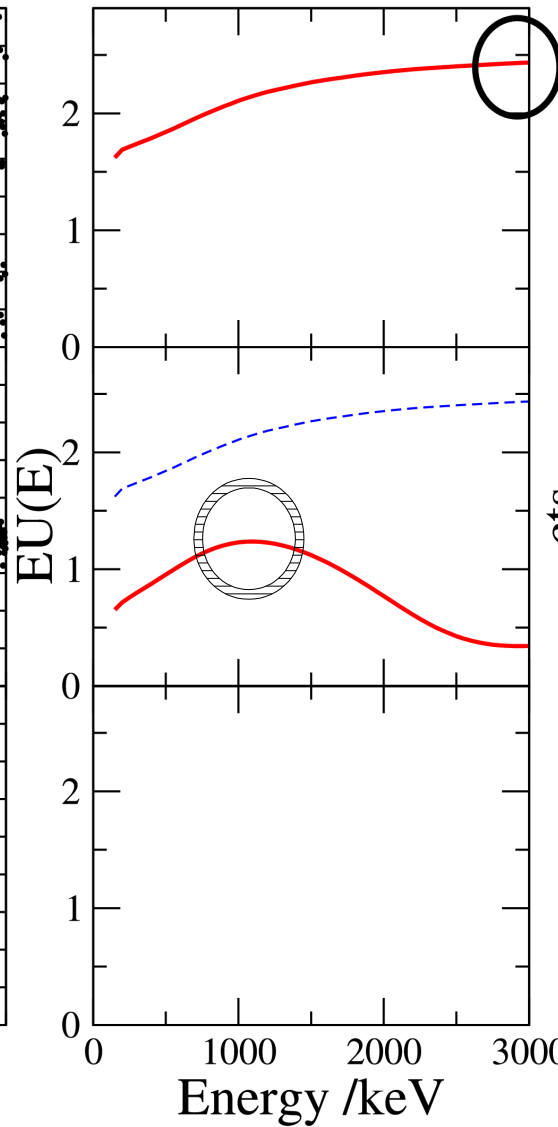
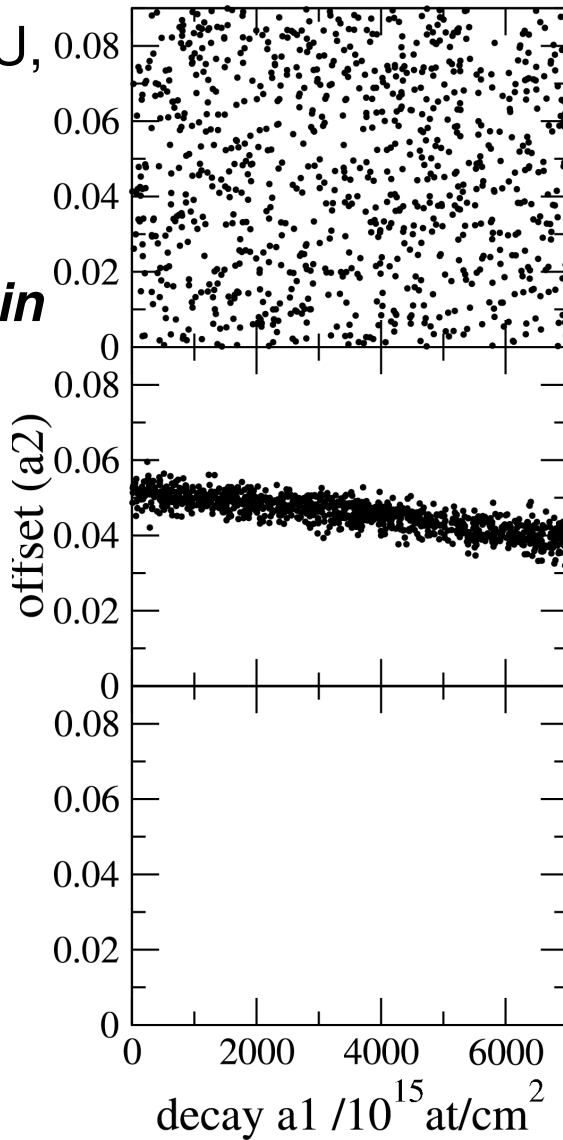
Compute EU(E),
and select best
energy

Integration by
posterior
sampling
CPU: 1-4 min



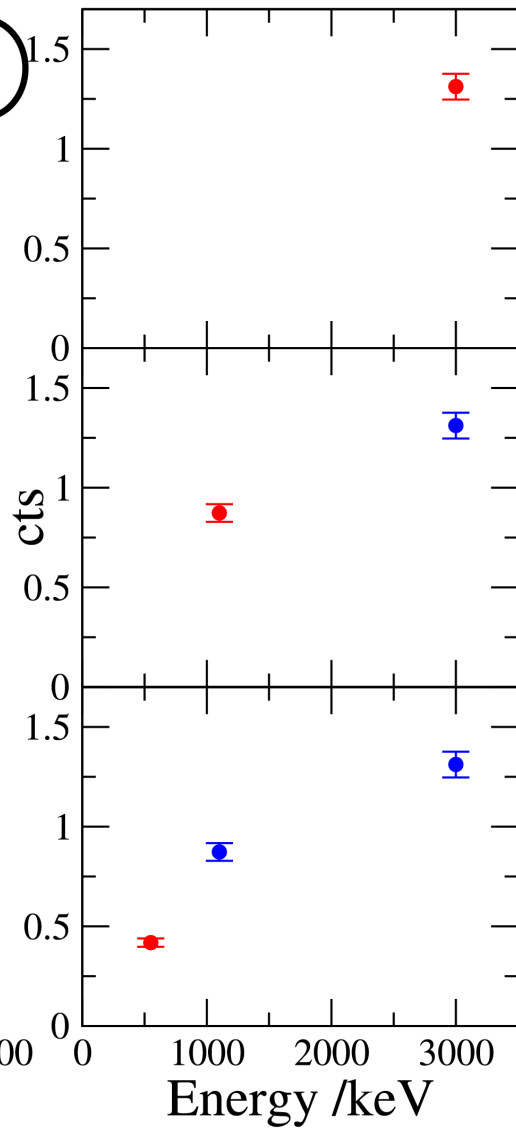
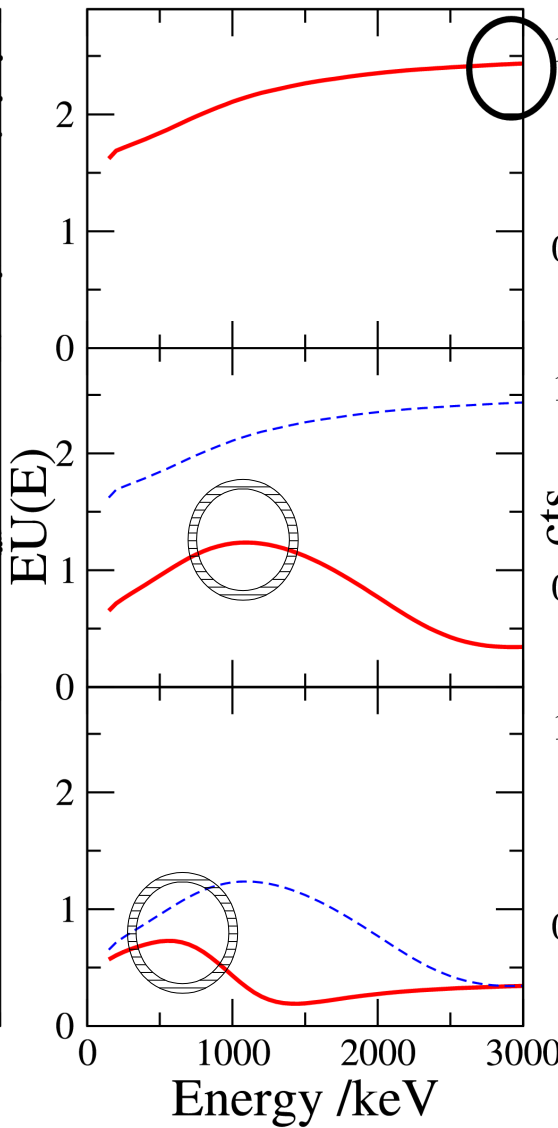
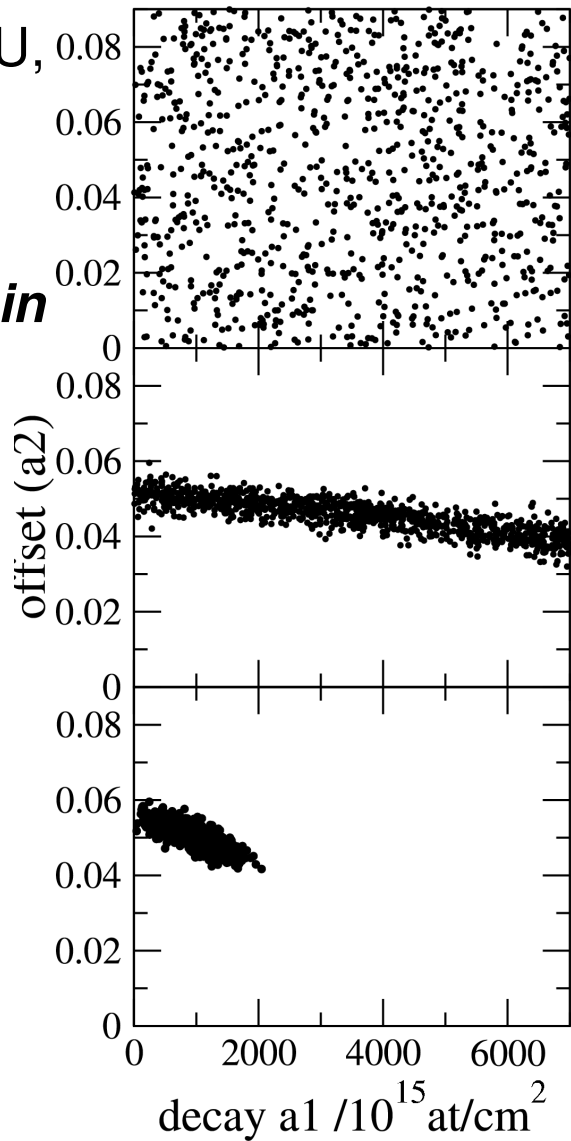
Optimize EU,
using
posterior
sampling
CPU:1-4 min

Cycle:
Prediction -
Verification



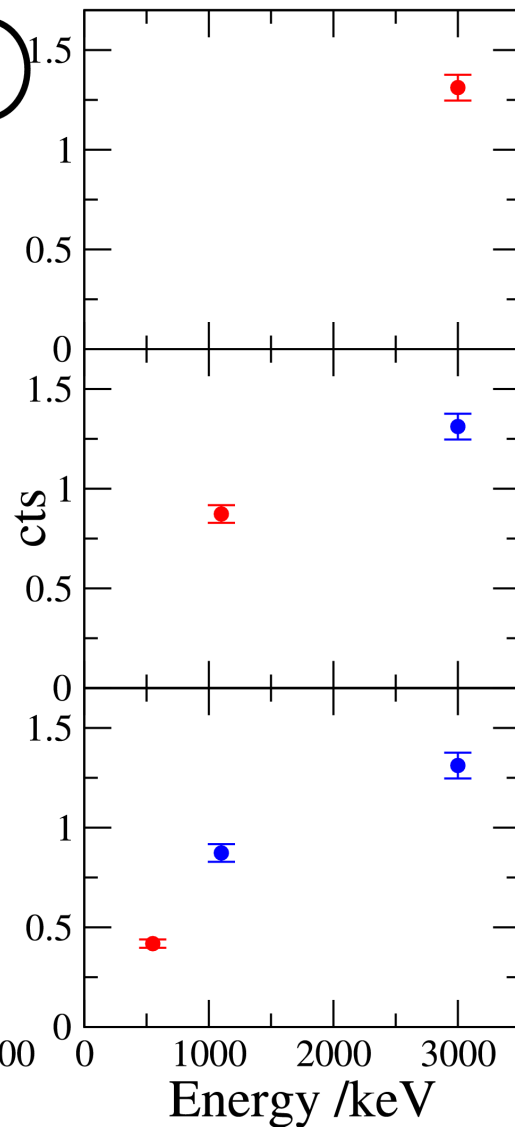
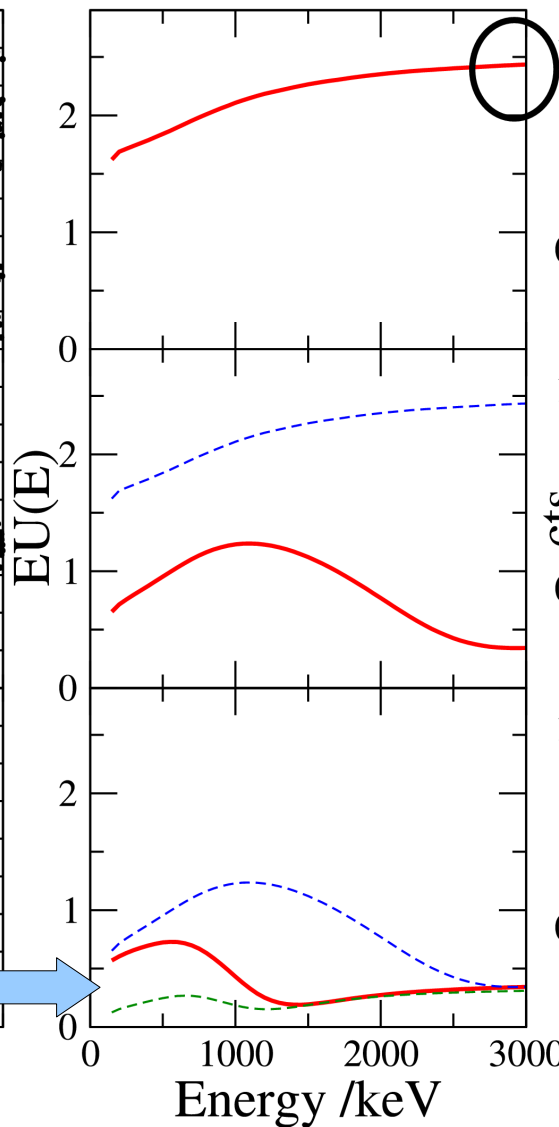
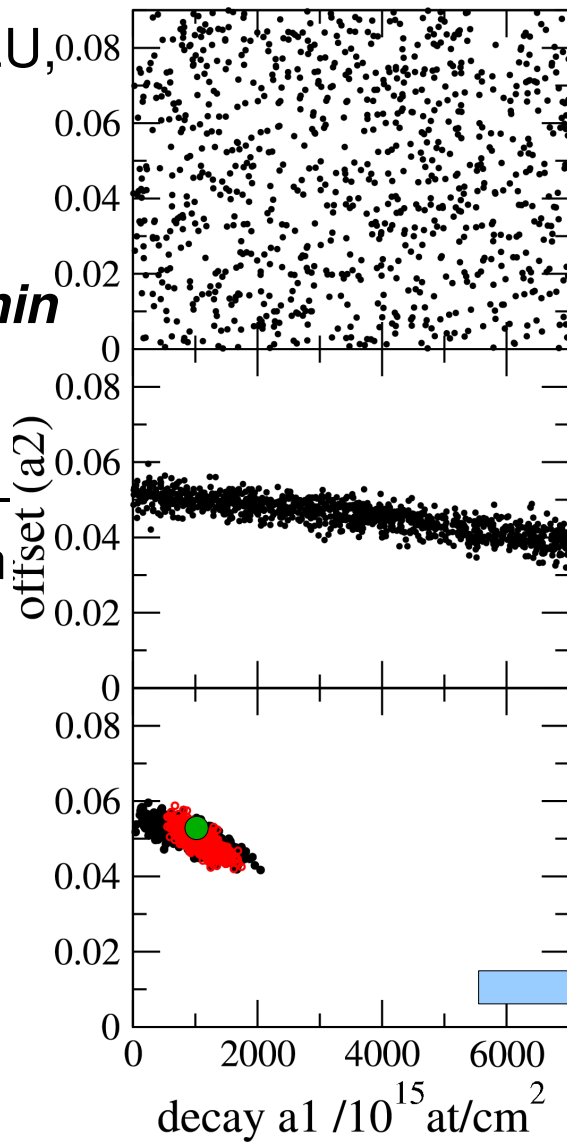
Optimize EU,
using
posterior
sampling
CPU:1-4 min

Cycle:
Prediction -
Verification



Optimize EU,
using
posterior
sampling
CPU:1-4 min

Cycle:
Prediction
Verification



Experimental Design


- Quantification of information

Value of measurements /diagnostics / experimental setups accessible
Optimized measurements protocols (on-line sequential design)

- Linear and **nonlinear** problems can be tackled

- Necessary integrations computationally demanding:

High-dimensional integrals in data- and parameter space
Anything beyond greedy algorithms unexplored (n-step ahead designs)

- Considerable gain in efficiency and/or accuracy has been demonstrated
- **Result-driven** automated measurement strategies (robotics) conceivable

IV. Numerical Interlude

Nested Sampling

Bayesian inference depends on (high-dimensional) **integration...**

$$\underbrace{p(D | \theta, \mathcal{H})}_{\text{likelihood}} \underbrace{p(\theta | \mathcal{H})}_{\text{prior}} = \underbrace{p(\theta | D, \mathcal{H})}_{\text{posterior}} \underbrace{p(D | \mathcal{H})}_{\text{evidence}}$$

posterior

$$p(\theta | D, \mathcal{H}) = \frac{p(D | \theta, \mathcal{H}) p(\theta | \mathcal{H})}{p(D | \mathcal{H})}$$

evidence

$$p(D | \mathcal{H}) = \int_{\theta} p(D | \theta, \mathcal{H}) p(\theta | \mathcal{H}) d\theta$$

model comparison

$$\begin{aligned} p(\mathcal{H}_1 | D) &\propto p(D | \mathcal{H}_1) p(\mathcal{H}_1) \\ p(\mathcal{H}_2 | D) &\propto p(D | \mathcal{H}_2) p(\mathcal{H}_2) \end{aligned}$$

Nested Sampling

Relation to statistical physics: partition function

Free energy

Probability of macrostate \mathcal{H} is proportional to:

$$Z(\beta, \mathcal{H}) = \int_{\theta} \exp(-\beta E(\theta, \mathcal{H})) d\theta$$

Evidence

How well a model \mathcal{H} predicted the data:

$$p(D | \mathcal{H}) = \int_{\theta} p(D | \theta, \mathcal{H}) p(\theta | \mathcal{H}) d\theta$$

Common problem:

$$Z = \int_{\theta} L(\theta) P(\theta) d\theta$$

$$L(\theta) \equiv \exp(-\beta E(\theta, \mathcal{H}))$$

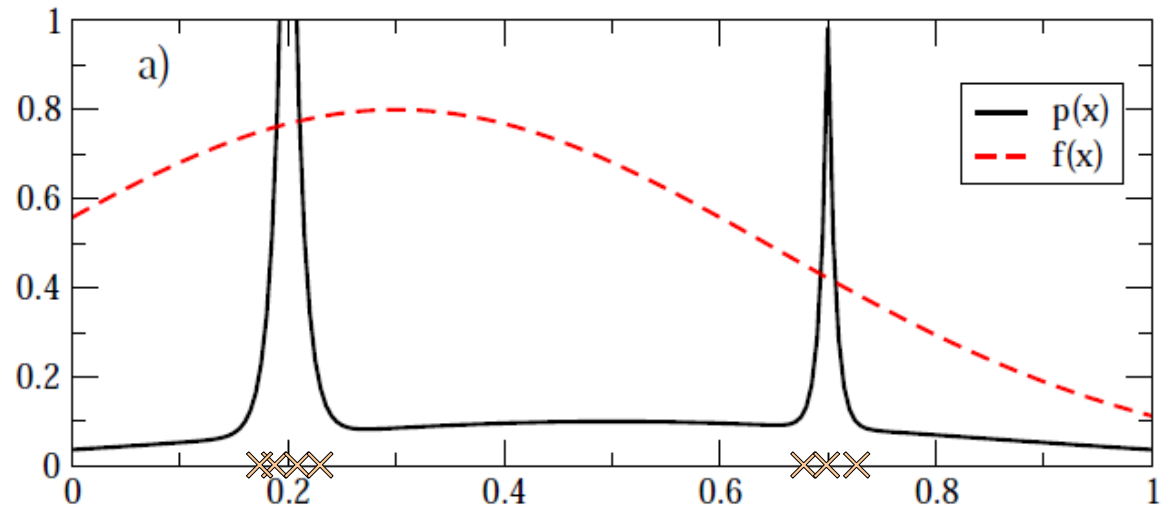
Boltzmann factor

$$L(\theta) \equiv p(D | \theta, \mathcal{H})$$

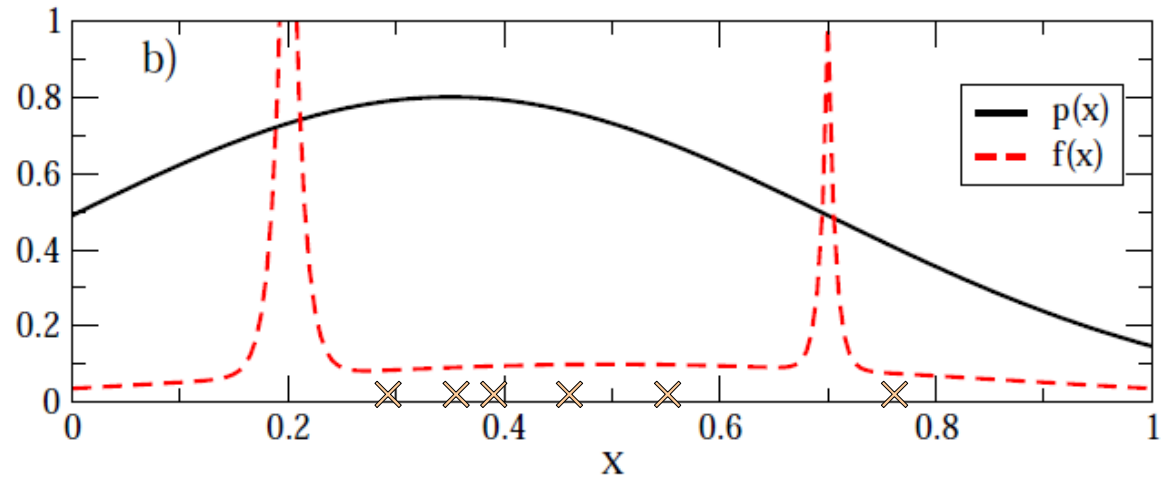
Likelihood function

Why is this integration difficult?

Typical situation
for expectation
values $\langle f \rangle$:



Typical situation
for evidence
calculation:



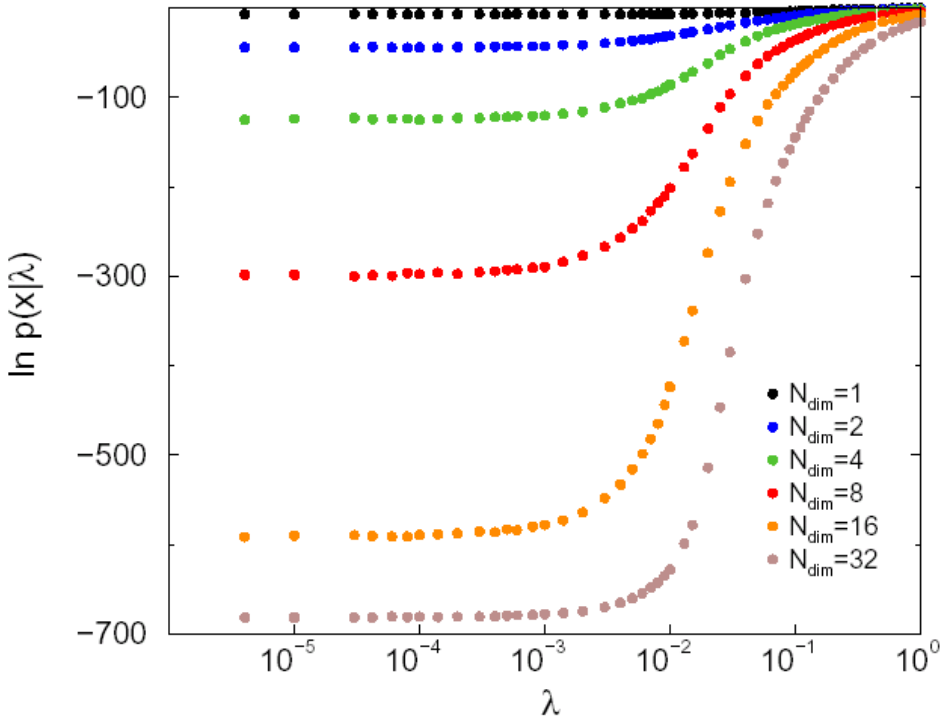
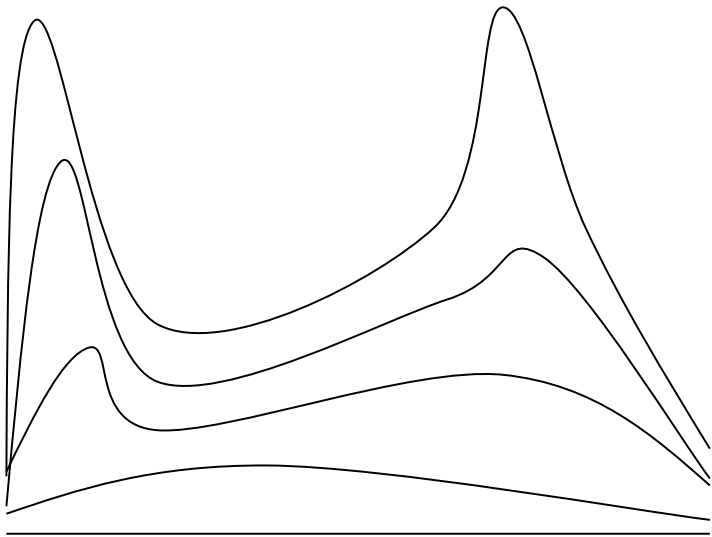
Nested Sampling

Thermodynamic Integration: Slowly introduce likelihood structure into integrand (similar: parallel tempering, perfect tempering):

$$Z(\lambda) = \int d\underline{x} \Gamma^\lambda(\underline{x}) \pi(\underline{x})$$

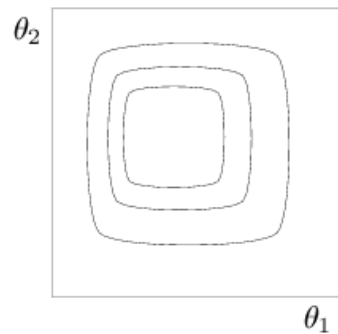
$Z(0)=1$ and $Z(1)=\text{Evidence}$

$$\ln(I) = \int_{\lambda=0}^{\lambda=1} d\lambda \frac{\partial \ln Z(\lambda)}{\partial \lambda} = \int_{\lambda=0}^{\lambda=1} d\lambda \int d\underline{x} \ln \Gamma(\underline{x}) \rho_\lambda(\underline{x})$$



R. Preuss, U. von Toussaint, AIP, 2007

$$Z = \int_{\theta} L(\theta) P(\theta) d\theta$$



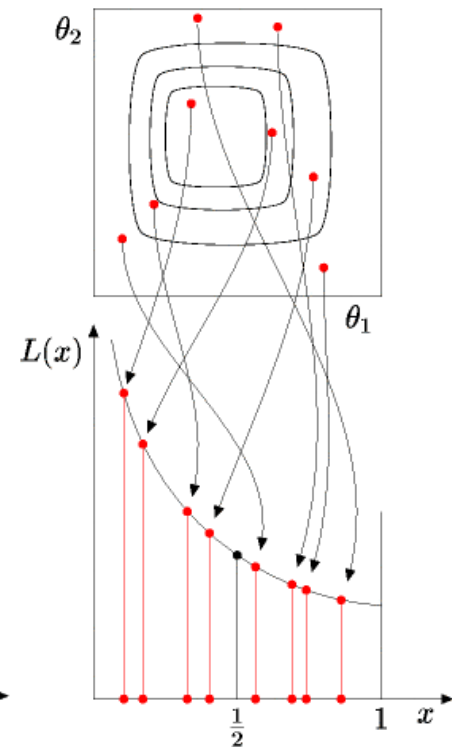
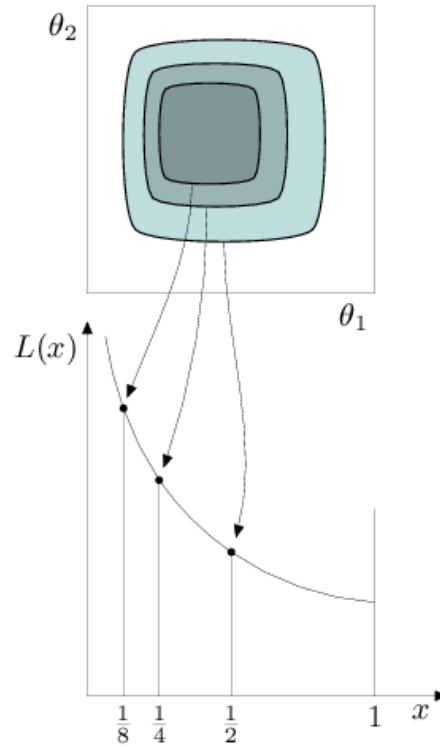
Contour plot of L

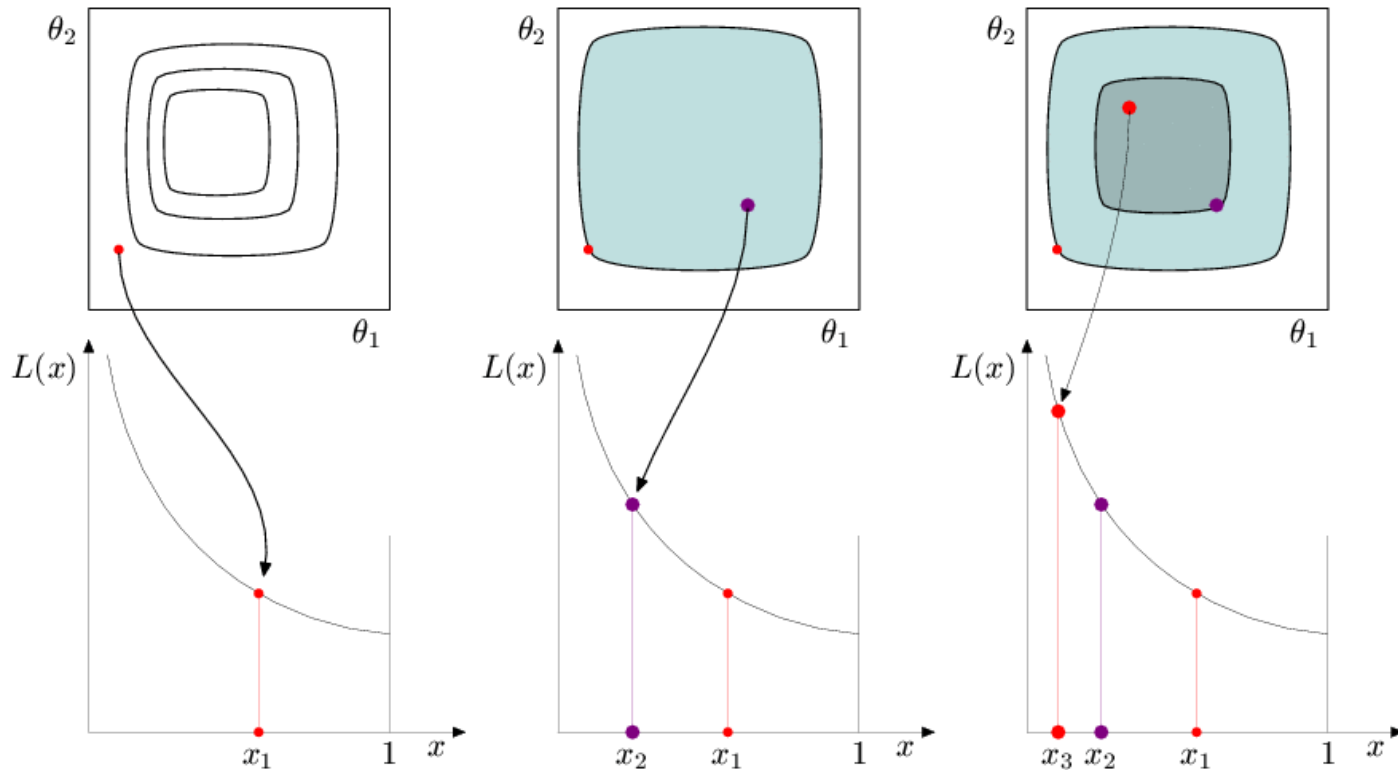
$$Z = \int_{\theta} L(\theta) P(\theta) d\theta = \int_0^1 L(x) dx$$

$$dx = P(\theta) d\theta$$

● Key concept:

● Sort all points by L





$$P(\theta^{(1)}) = P(\theta)$$

$$P(\theta^{(i+1)}) \propto \begin{cases} P(\theta) & L(\theta) > L_i \\ 0 & \text{otherwise} \end{cases}$$

$$x_1 \sim \text{Uniform}(0, 1)$$

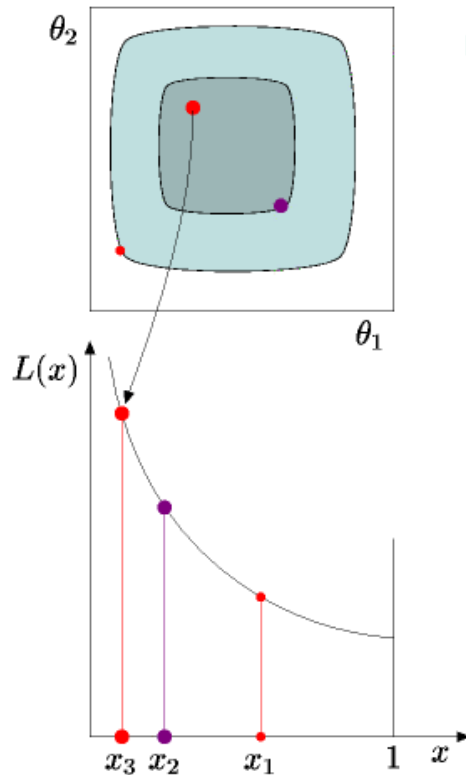
$$x_2 \sim \text{Uniform}(0, x_1)$$

$$x_3 \sim \text{Uniform}(0, x_2)$$

$$\langle x_1 \rangle = \frac{1}{2}$$

$$\langle x_2 \rangle = \frac{1}{2} \langle x_1 \rangle = \frac{1}{4}$$

$$\langle x_3 \rangle = \frac{1}{8}$$



● Draw from

$$P(\theta^{(i+1)}) \propto \begin{cases} P(\theta) & L(\theta) > L_i \\ 0 & \text{otherwise} \end{cases}$$

- cf Annealing's intermediate distributions

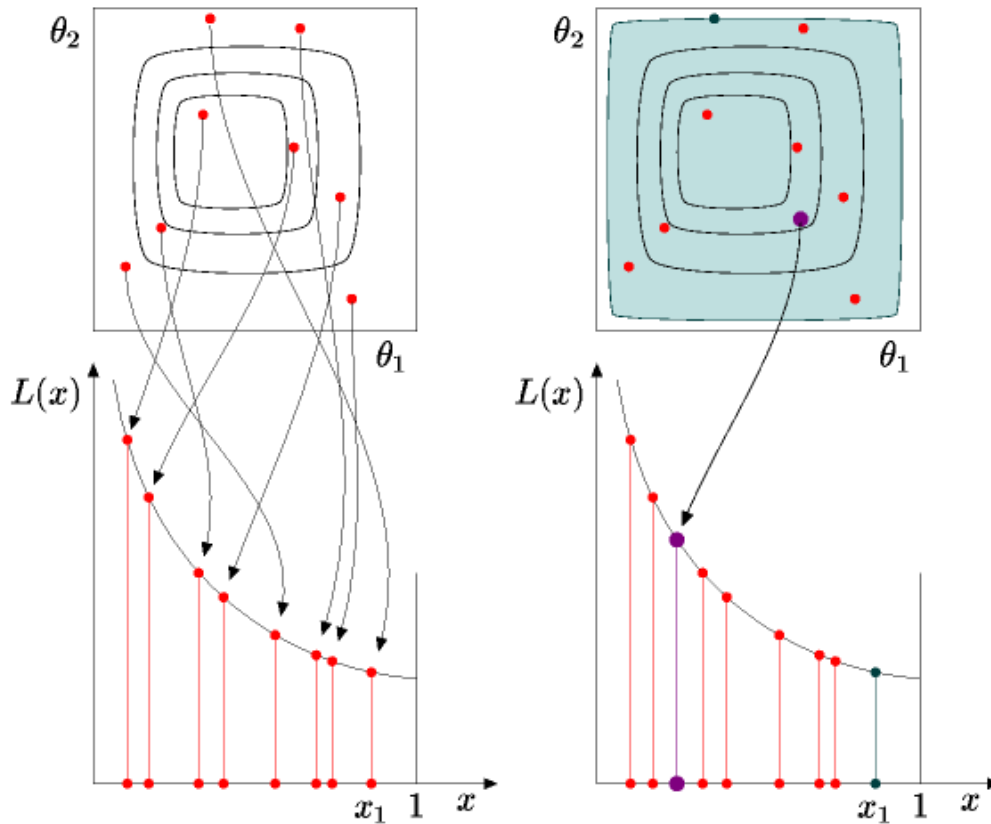
$$P(\theta | \beta) \equiv \frac{1}{Z(\beta)} L(\theta)^\beta P(\theta)$$

- Assuming uniform sampling subject to an energy constraint is possible, $P(\mathbf{X})$ is simple

$$\left\langle \log \frac{x_{i+1}}{x_i} \right\rangle = -1 \quad \text{Order statistics for } P(x)$$

independent of L

- Nested sampling's behaviour is invariant under monotonic transformations of L



The result:

$$\hat{\mathcal{Z}} = \sum_i \delta \hat{x}_i L_i$$

$$\hat{x}_i \equiv \exp(-i/N)$$

Pro: Quite different approach, very general, easy to implement

Con: Uniform sampling under constraint?

V. Conclusion

Summary
Outlook

- Bayesian probability theory provides consistent and transparent approach to the cycle of scientific inference
- Incorporation of available prior information is straightforward
- Drawback of numerical complexity mitigated by
 - New algorithms
 - Increasing computing power
- State of the Art: parameter estimation, model comparison
- Coming soon: probabilistic combination of diagnostics (IDA)
- Still (largely) unexplored:
 - potential of Experimental Design, e.g. robotics, self-adapting 'measurement'-strategies on computer simulations (e.g. automated MD potential generation)
 - novelty detection in large scale simulations / experiments

- Promising and/or unexplored research directions

Bayesian Experimental Design

Large data sets and complex (simulation based models) require consistent response-surface estimates (O'Hagan)

Estimation of stochastic partial differential equations or functionals from data (e.g. for turbulence)

Data based model design and model estimation, e.g. with respect to possible causation instead of correlation only (Pearl)

- 32th Workshop on Bayesian Inference and Maximum Entropy Methods at IPP Garching (15.-20. July 2012)
See: <http://www.ipp.mpg.de/maxent2012>

32th Workshop
on
Bayesian Inference and Maximum Entropy Methods
at
IPP Garching (15.-20. July 2012)

See: <http://www.ipp.mpg.de/maxent2012>

Thank you
for
your attention