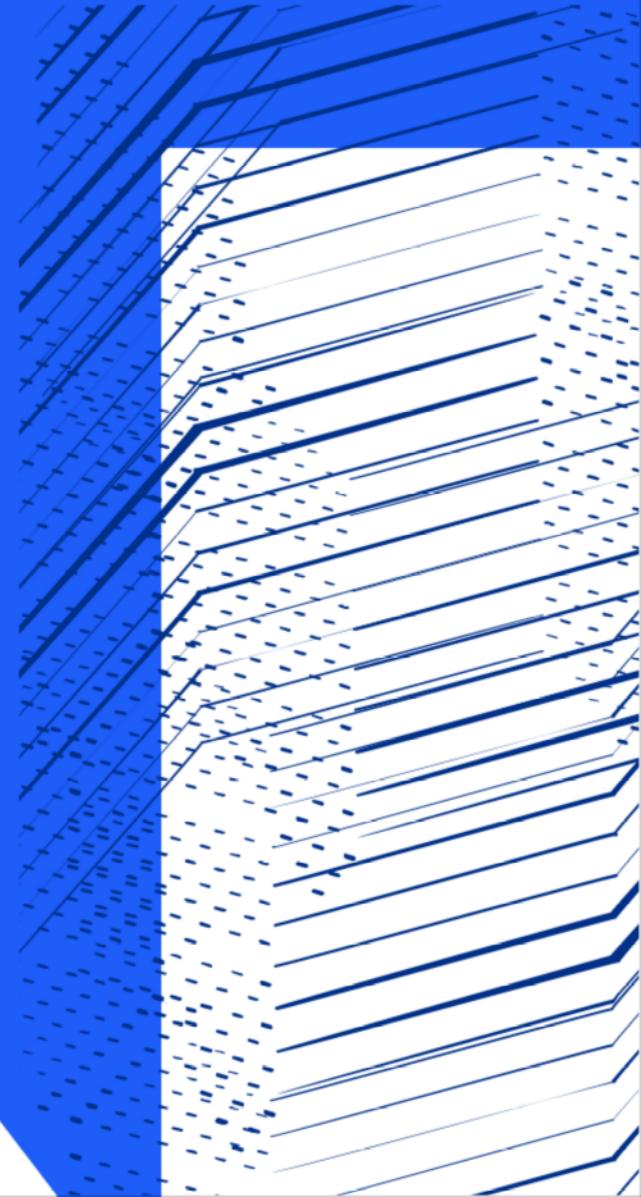




Science and
Technology
Facilities Council

Stochastic Model Fitting

Taking methods back to physical basics



Science and
Technology
Facilities Council

Agenda

1 Goals

The what and why.

2 Methods

3 Procedures

Ensuring and documenting progress



Science and
Technology
Facilities Council

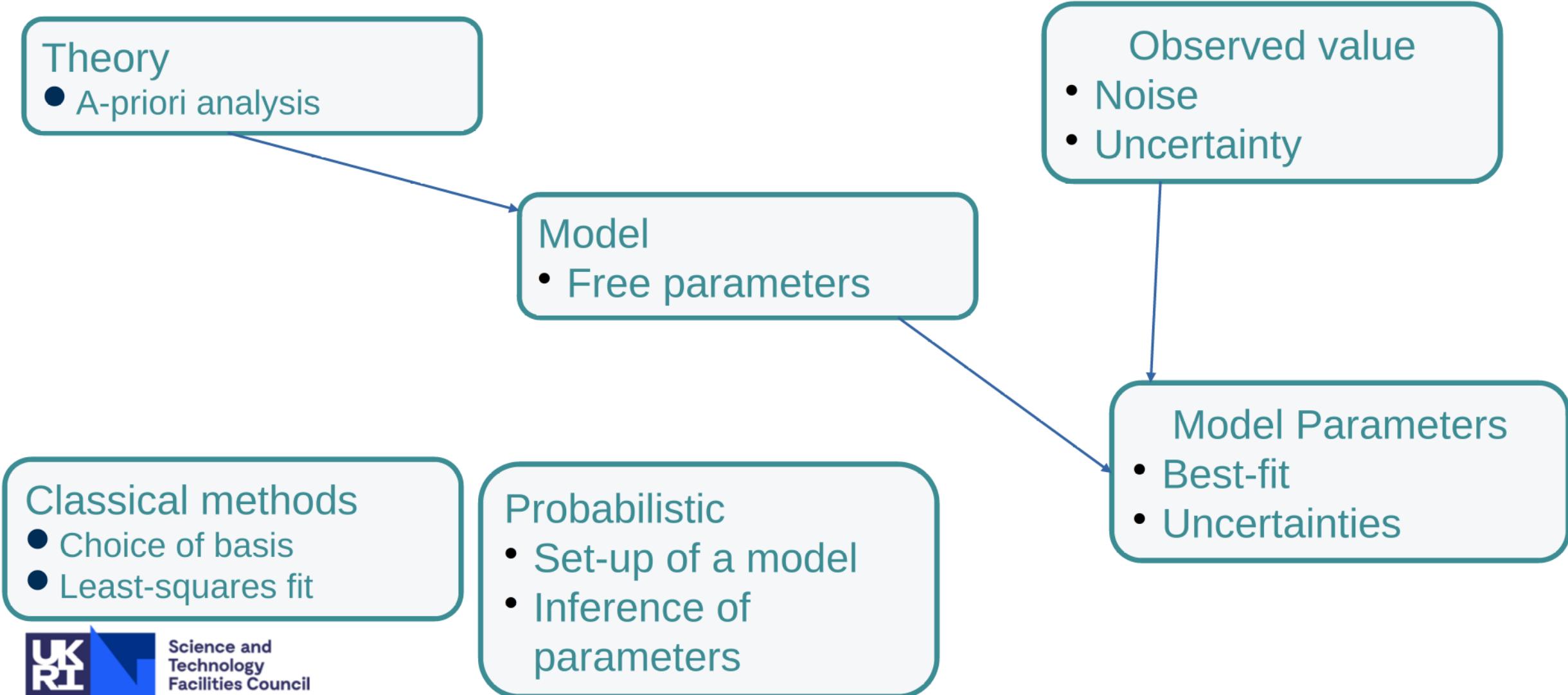


Image © STFC Alan Ford

The Overarching Goal

What use case do we consider?

Fitting probabilistic models to data



Inference from measurements – Bayes' rule

Probability of B occurring if A has occurred
– given a model

The same, but with densities

$$P[B_i|A] = \frac{P[A|B_i]P[B_i]}{\sum_{j=1}^n P[A|B_j]P[B_j]}$$

$$p_{X|Y}(x|y) = \frac{p_X(x)p_{Y|X}(y|x)}{p_Y(y)}$$

Not solvable by hand for a n arrow's flight

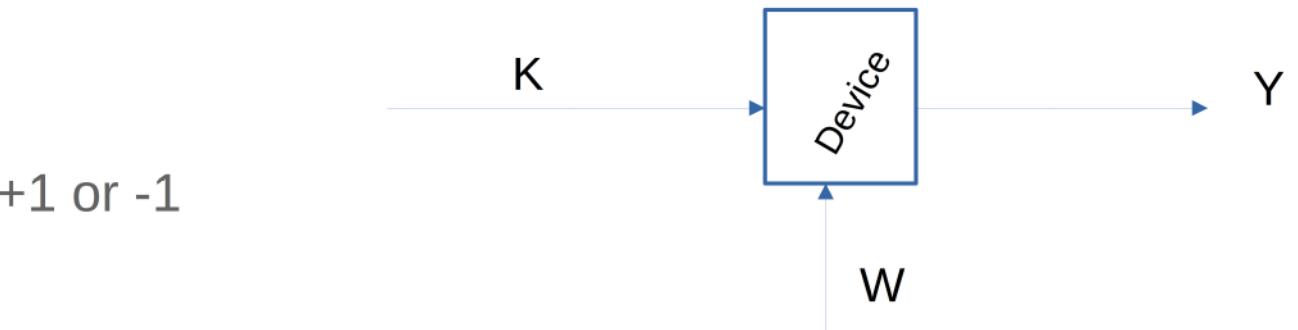
$$Z_i = Z_0 + V_0(t - t_0) + A_0(t - t_0)^2 + W_i$$
$$p_{(Z,V,A)|(z_1, \dots, z_n)}((z, v, a)|(t_1, \dots, t_n)) = \frac{\prod_{i=1}^3 p_{(Z,V,A)_i}((z, v, a)_i) \prod_{j=1}^N p_{Z_i|(Z,V,A)}(z_i|(z, v, a))}{p_{(Z,V,A)}((z, v, a))}$$

Bayes' rule example: Noisy digital output

Discrete signal with noise

- Input: signal value K equally likely to be $+1$ or -1
- Additive, normally distributed noise W
- Measure $Y = K + W$
- Given measurement y what is the probability of the signal K being 1 ?

$$p_K(k) = \frac{1}{2}$$
$$p_{Y|K}(y|k) = \frac{1}{\sqrt{2\pi}} e^{\frac{-1}{2}(y-k)^2}$$



$$\rightarrow p_Y(y) = \frac{1}{2\sqrt{2\pi}} e^{\frac{-1}{2}(y+1)^2} + \frac{1}{2\sqrt{2\pi}} e^{\frac{-1}{2}(y-1)^2}$$

$$p_{K|Y}(1|y) = \frac{1}{1+e^{-2y}}$$



Science and
Technology
Facilities Council

Methods

What will we be developing?



Science and
Technology
Facilities Council

Methods to fit a model to existing data

Deterministic

- Polynomial fitting
- Polynomial interpolation
- Trigonometric Interpolation

Statistics

- Machine Learning
- Bayesian Inference
- Markov Chain methods

Determining the minimum/ground state of an unknown potential



Problems of the common methods

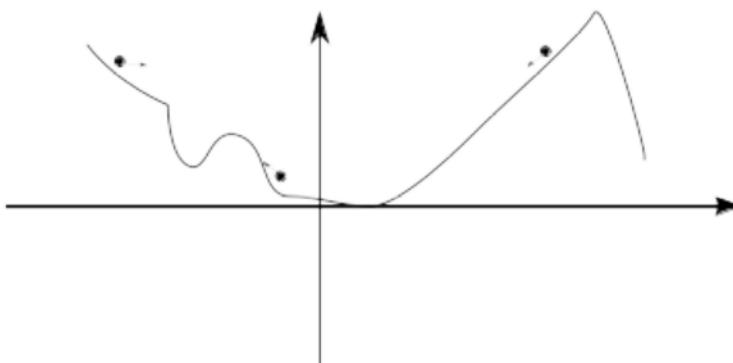
Deterministic

- Overfitting
- Numerical stability
- Computational cost

Statistics

- Large datasets
- Slow convergence
- Computational cost

Solution: Consider the problem as an energy minimization problem
Finding the minimum of an arbitrary potential.



Solving parameter estimation by finding the ground state of a Hamiltonian

$$H(q, p) = T(p) + V(q) = \sum p_i^2 / (2m_i) + V(q_i)$$

Random Walk

- Direction and step size chosen at random
- Step (almost) always accepted

Markov Chain MC

- Fixed step size
- Direction chosen at random
- Metropolis-Hastings for acceptance

Hamiltonian MC

- Sampling trajectories in phase space
- Extensive propagation
- Metropolis-Hastings / Gibbs



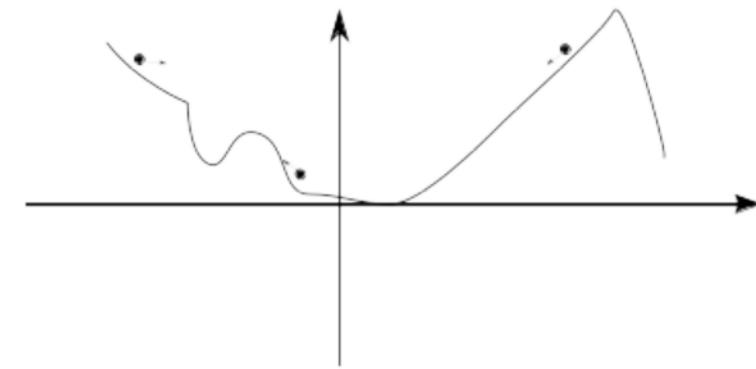
Science and
Technology
Facilities Council

Hamiltonian Monte Carlo – the method of choice

$$H(q, p) = T(p) + V(q) = \sum p_i^2 / (2m_i) + V(q_i)$$

Algorithm 1 Hamiltonian Monte Carlo

```
Given  $\theta^0, \epsilon, L, \mathcal{L}, M$ :  
for  $m = 1$  to  $M$  do  
    Sample  $r^0 \sim \mathcal{N}(0, I)$ .  
    Set  $\theta^m \leftarrow \theta^{m-1}, \bar{\theta} \leftarrow \theta^{m-1}, \bar{r} \leftarrow r^0$ .  
    for  $i = 1$  to  $L$  do  
        Set  $\bar{\theta}, \bar{r} \leftarrow \text{Leapfrog}(\bar{\theta}, \bar{r}, \epsilon)$ .  
    end for  
    With probability  $\alpha = \min \left\{ 1, \frac{\exp\{\mathcal{L}(\bar{\theta}) - \frac{1}{2}\bar{r} \cdot \bar{r}\}}{\exp\{\mathcal{L}(\theta^{m-1}) - \frac{1}{2}r^0 \cdot r^0\}} \right\}$ , set  $\theta^m \leftarrow \bar{\theta}, r^m \leftarrow -\bar{r}$ .  
end for  
  
function Leapfrog( $\theta, r, \epsilon$ )  
    Set  $\bar{r} \leftarrow r + (\epsilon/2)\nabla_\theta \mathcal{L}(\theta)$ .  
    Set  $\bar{\theta} \leftarrow \theta + \epsilon\bar{r}$ .  
    Set  $\bar{r} \leftarrow \bar{r} + (\epsilon/2)\nabla_\theta \mathcal{L}(\bar{\theta})$ .  
    return  $\bar{\theta}, \bar{r}$ .
```

$$P_C[(q, p)] = \frac{1}{Z} e^{-\beta H(q, p)}$$
$$Z = \sum e^{-\beta H(q_i, p_i)}$$


Hamiltonian Monte Carlo – problems

$$H(q, p) = T(p) + V(q) = \sum p_i^2 / (2m_i) + V(q_i)$$

Algorithm 1 Hamiltonian Monte Carlo

```
Given  $\theta^0, \epsilon, L, \mathcal{L}, M$ :  
for  $m = 1$  to  $M$  do  
    Sample  $r^0 \sim \mathcal{N}(0, I)$ .  
    Set  $\theta^m \leftarrow \theta^{m-1}, \bar{\theta} \leftarrow \theta^{m-1}, \bar{r} \leftarrow r^0$ .  
    for  $i = 1$  to  $L$  do  
        Set  $\bar{\theta}, \bar{r} \leftarrow \text{Leapfrog}(\bar{\theta}, \bar{r}, \epsilon)$ .  
    end for  
    With probability  $\alpha = \min \left\{ 1, \frac{\exp\{\mathcal{L}(\bar{\theta}) - \frac{1}{2}\bar{r} \cdot \bar{r}\}}{\exp\{\mathcal{L}(\theta^{m-1}) - \frac{1}{2}r^0 \cdot r^0\}} \right\}$ , set  $\theta^m \leftarrow \bar{\theta}, r^m \leftarrow -\bar{r}$ .  
end for  
  
function Leapfrog( $\theta, r, \epsilon$ )  
    Set  $\bar{r} \leftarrow r + (\epsilon/2)\nabla_\theta \mathcal{L}(\theta)$ .  
    Set  $\bar{\theta} \leftarrow \theta + \epsilon\bar{r}$ .  
    Set  $\bar{r} \leftarrow \bar{r} + (\epsilon/2)\nabla_\theta \mathcal{L}(\bar{\theta})$ .  
    return  $\bar{\theta}, \bar{r}$ .  
  
 $P_C[(q, p)] = \frac{1}{Z} e^{-\beta H(q, p)}$   
 $Z = \sum e^{-\beta H(q_i, p_i)}$ 
```

Methodological

- Step size
- Number of steps
- Initial conditions (position, momentum)

Pedagogical

- Mathematical mutilation

Computational

- Stability dictated by integrator
- Expensive evaluation of potentials
- Parameter ranging
- Mediocre scalability
- Synchronisation for global averages

Ensemble-based extension – and other improvements

$$H(q, p) = T(p) + V(q) = \sum p_i^2 / (2m_i) + V(q_i)$$

Ensemble-based HMC

- μ -canonical ensemble
- Ensemble-averages
- Divide-and-conquer

Extensions

- GPU offloading
- Variational Quantum Encoding
- Simlated/Quantum Annealing

Physical extensions

- Canonical ensemble
- Variable step-size
- Optimal step-size choice
- Initialization
- Max-entropy principle



Science and
Technology
Facilities Council



Science and
Technology
Facilities Council

Procedures

The technical part.



Science and
Technology
Facilities Council

Meetings

- **Times:**
 - **Monday** 0900BST / 1000 CEST – for 1 h
 - **Thursday** 1400 BST / 1500 CEST – for 30 min – 1h
 - Upon request.
- **Mode**
 - Video conference via Zoom.

Basic tools – GIT, Python

Programming:

- Python (possibly C++)
- NumPy/SciPy
- **NumPyro**
(<https://num.pyro.ai/en/stable/index.html>)



GIT rules:

- **No binaries**
- **Pull-Request only**
- **No acceptance without a test**
- Questions as issues.
- At least one approval required before merge.
- Method summaries as TeX files.

Collaboration / communication

- GIT repository
- Email + VC

Theoretical basics – literature.

Physics:

- Statistical physics and thermodynamics
 - Lecture notes by M. Sigrist (ETH Zürich Fall term 2012)
 - Elementary Statistical Physics (C. Kittel, Dover publications)
- Classical mechanics
 - Lecture notes on Classical Dynamics (D. Tong, U Cambridge)
 - Theoretical Physics I (L.D. Landau)
 - Fundamentals of Mechanics (J. Marsden)

Mathematics:

- Mathematical Foundations of Statistical Mechanics (A. Khinchin, Dover publications)
- „Efficient Sequential Monte-Carlo
- Samplers for Bayesian Inference“ (Thi Nguyen et al., IEEE transactions on signal processing 64, 2016)
- „The No-U-Turn sampler“ (Hoffman et al, Journal of Machine Learning Research 15, 2014)

Methods:

- Monte Carlo Simulation in Statistical Physics (Binder, Heerman, Springer GTP)



Science and
Technology
Facilities Council



A large, bold, white question mark is centered on a blue rectangular background. The background is surrounded by a series of blue arrows pointing towards the right, creating a sense of motion and direction.

Questions?



Science and
Technology
Facilities Council