

Imperial-X

Machine Learning Basics for Regression and Inference

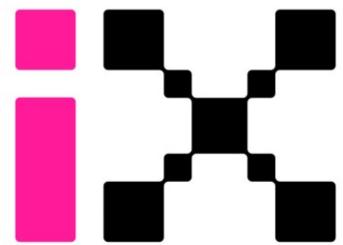


Code and slides
URL:
[https://github.com/aidancrilly
/ML_Lecture_Demos](https://github.com/aidancrilly/ML_Lecture_Demos)

Opinion on machine learning in your research:

Go to menti.com and enter code 87 90 96 11

0	0	0	0	0
What even is ML?	Not useful	Unsure of utility	Would like to use in future	Actively using it



Imperial-X

Machine Learning Basics for Regression and Inference

Dr Aidan Crilly

Schmidt AI in Science Postdoctoral Fellow

What to expect?

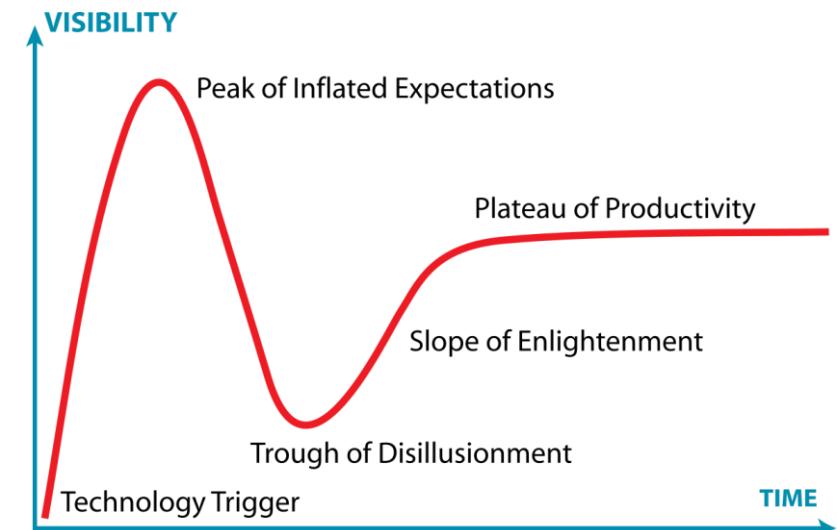
- This talk aims to:
 - Introduce key machine learning (ML) concepts and language
 - Relate ML techniques to more familiar numerical methods
 - Introduce ML techniques which have been used in Plasma Group's research
 - Hopefully, be a jumping off point to learning more about ML useful to your own research
- This talk does not aim to:
 - Discuss ML research, we will discuss ML *in* research
 - Be an authority on ML, it would be great if this could become a discussion on use cases

Overview

- What is Machine Learning?
- Regression & Classification
- Parametric Regressors (Non-neural)
 - Ordinary and Non-Linear Least Squares
- Bayesian Inference
- Non-Parametric Regressors
 - Gaussian Processes
- Neural Networks

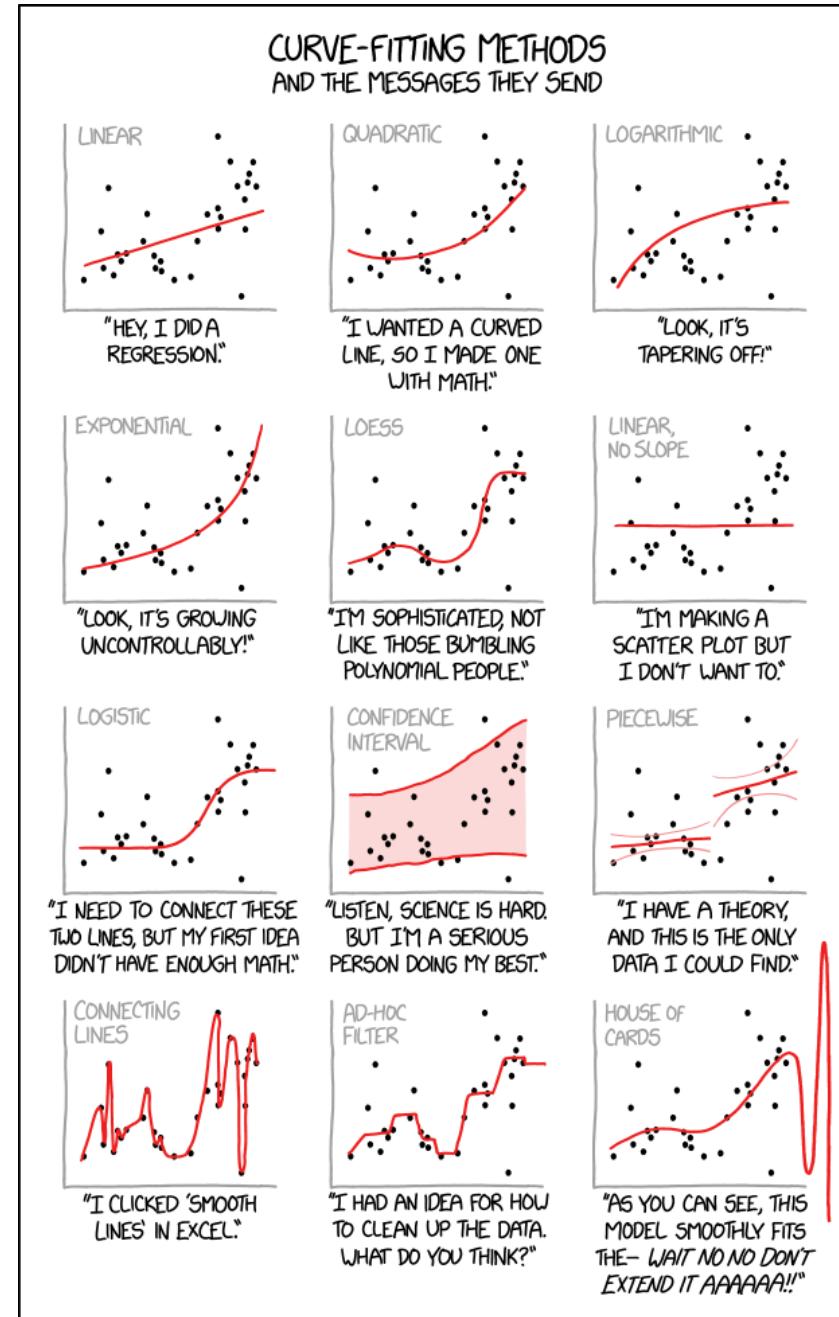
What is “Machine Learning”?

- Numerical modelling which adapts parameters algorithmically to learn the relationship between input and output data
- Parameter evolution often posed as an optimisation problem (nothing new!)
- User often must select ‘hyper-parameters’ which change aspects of the machine learning algorithm



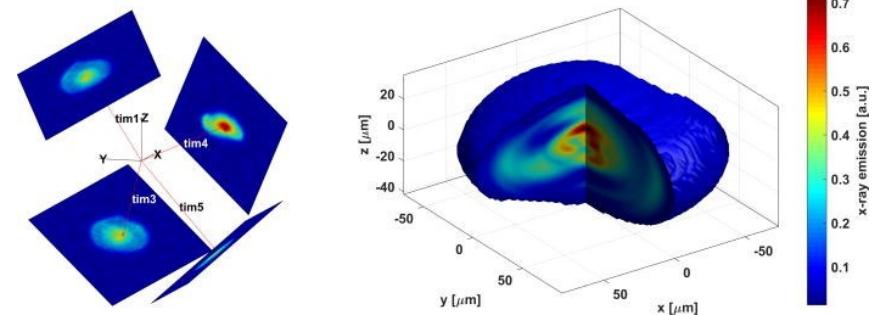
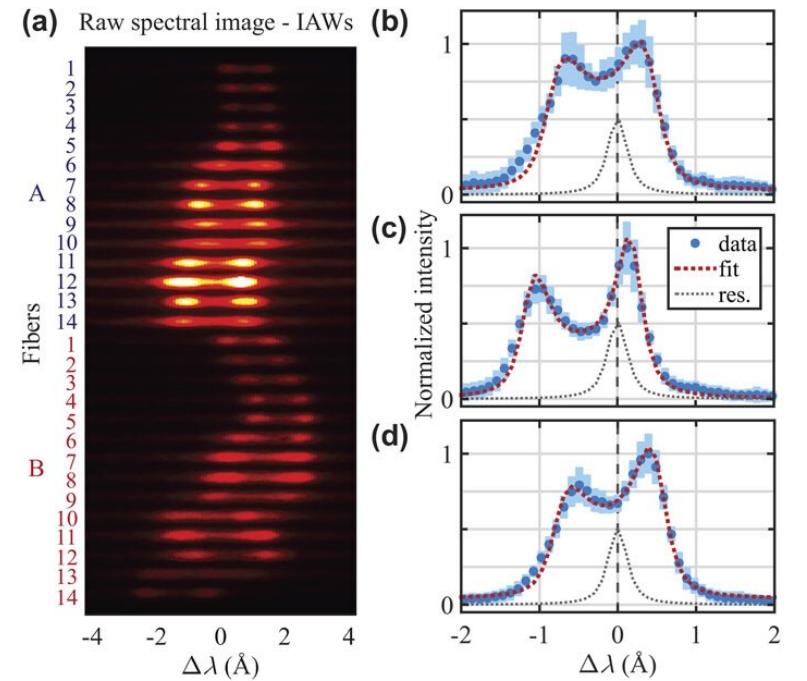
Regression & Classification

- Two of the core ML tasks are regression and classification
- Classification is the task of placing data within a finite number of sets, for example determining if a photo contains a cat or not
- Regression is the task of estimating the functional relationship between a set of input and output variables, for example fitting a straight line to (x,y) data pairs



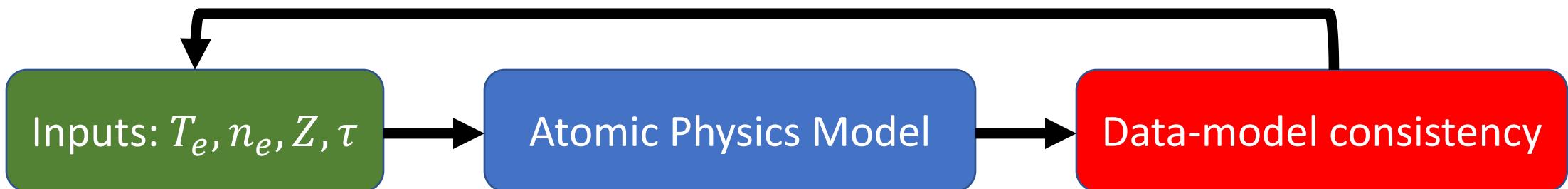
Regression

- Regression problems are ubiquitous in data analysis within the physical sciences
- An entirely non-exhaustive list of examples from my research:
 - Image reconstruction and 3D tomography
 - Thomson scattering spectra
 - Fusion neutron spectroscopy
 - X-ray line spectroscopy
- Typically, reduced models with small number of parameters are used to fit each diagnostic signal separately
- Latent parameters of numerical models can also be learned via regression



Parametric Regressors

- The most common form of regression includes a ‘physics-based’ model which maps input physical parameters (e.g. temperature and density) to the expected diagnostic signature
- For example, line spectroscopy



- The ‘optimal’ set of parameters which match the observed data are sought in a forward-fit process

Ordinary Least Squares



- Say we have a linear model with unknown parameters θ_j :

$$f(x, \theta) = \sum_j \theta_j f_j(x) = \underline{\mathbf{A}} \cdot \boldsymbol{\theta}$$

- And we wish to minimize the weighted least squares distance to the data y_i (we shall see why later):

$$\begin{aligned} \min_{\boldsymbol{\theta}} \sum_i w_i [y_i - f(x_i, \boldsymbol{\theta})]^2 &= \min_{\boldsymbol{\theta}} (\mathbf{y}^T \underline{\mathbf{W}} \mathbf{y} + \boldsymbol{\theta}^T \underline{\mathbf{A}}^T \underline{\mathbf{W}} \underline{\mathbf{A}} \boldsymbol{\theta} - 2 \boldsymbol{\theta}^T \underline{\mathbf{A}}^T \underline{\mathbf{W}} \mathbf{y}) \\ &\rightarrow \underline{\mathbf{A}}^T \underline{\mathbf{W}} \underline{\mathbf{A}} \hat{\boldsymbol{\theta}} = \underline{\mathbf{A}}^T \underline{\mathbf{W}} \mathbf{y} \end{aligned}$$

- Therefore, the ‘best fit’ is given by the solution to the ‘normal equations’:

$$\hat{\boldsymbol{\theta}} = (\underline{\mathbf{A}}^T \underline{\mathbf{W}} \underline{\mathbf{A}})^{-1} \underline{\mathbf{A}}^T \underline{\mathbf{W}} \mathbf{y}$$

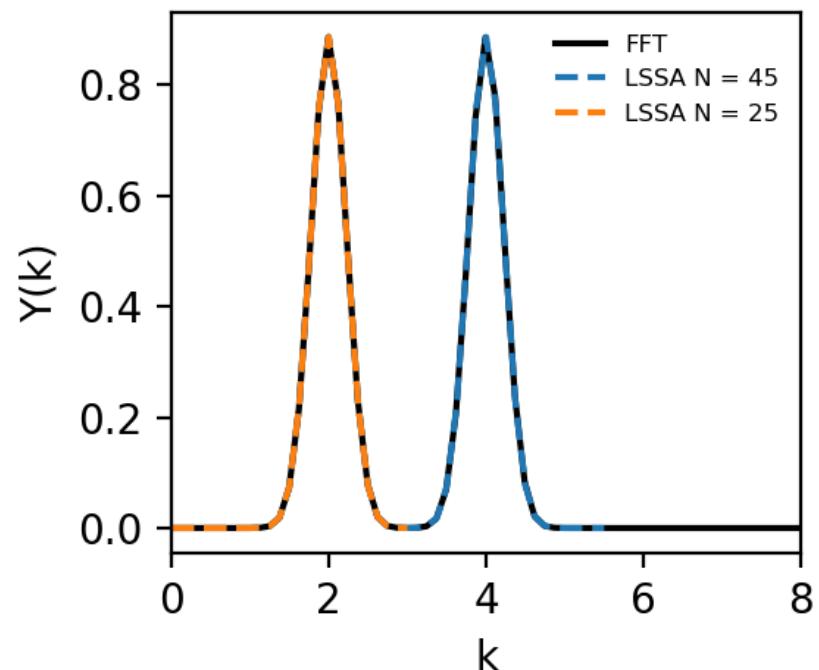
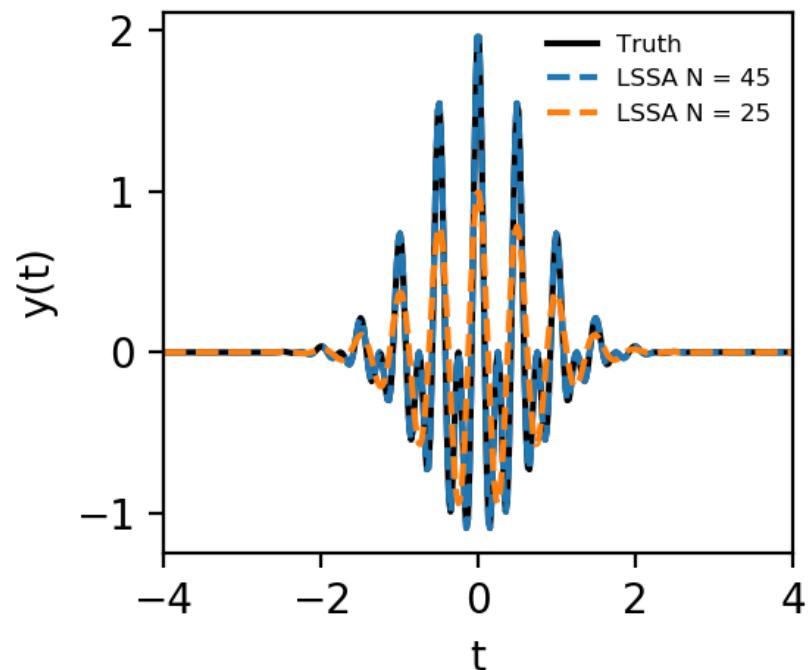
Practical example

- Least-squares spectral analysis (LSSA) is a linear regression problem

$$f(x, \theta) = \sum_{j=0}^N \theta_j \cos\left[\frac{\pi j}{T} t\right] = \underline{A} \cdot \underline{\theta}$$

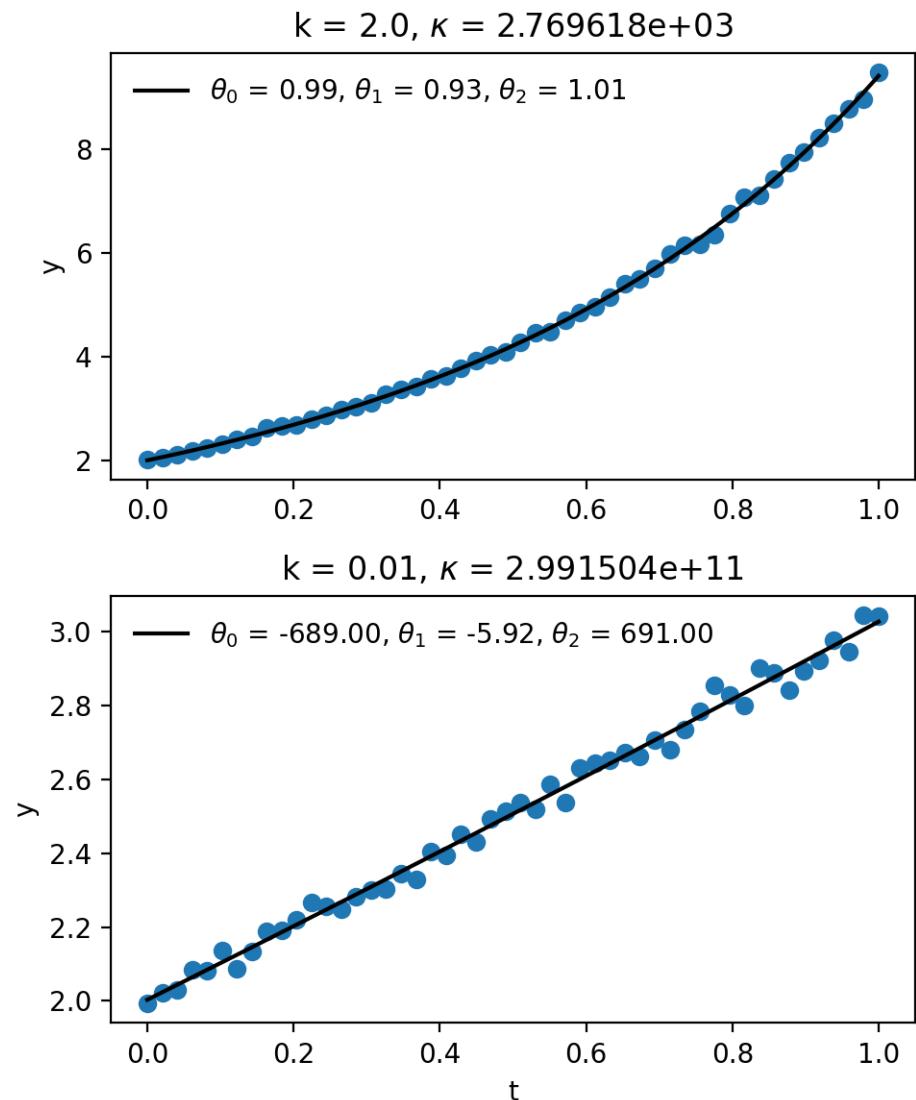
$$\hat{\underline{\theta}} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{y}$$

- Can handle unequally spaced data, uncertainties and limit the mode number
 - See Lomb-Scargle method
 - Not formally a Discrete Fourier Transform



Ill-conditioned or ill-posed problems

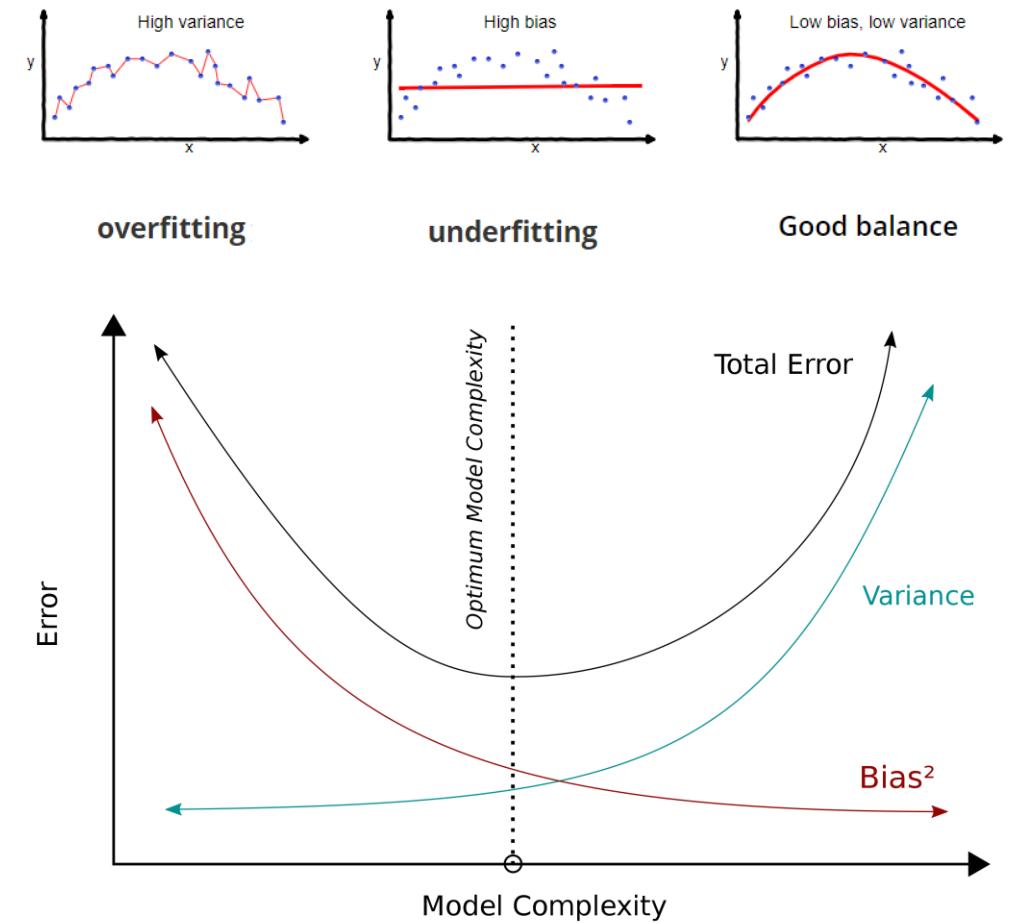
- While linear systems have formal solutions they can be numerical unstable, small errors on one side of the equation can lead to large errors on the other
- One possible cause is degeneracy or collinearity of input parameters, e.g. fitting $f(\theta, k, t) = \theta_0 + \theta_1 t + \theta_2 \exp[kt]$, for small kt
- These problems are ‘ill-conditioned’ as quantified by the condition number of the linear system



Regularisation – Bias-variance tradeoff

- Another key ML technique is regularisation
- It aims to reduce variance while introducing bias i.e. prevents overfitting
- Error on *unseen* sample:

$$\text{MSE} = \text{Var}(f) + \text{Bias}(f)^2 + \sigma^2$$



Regularisation

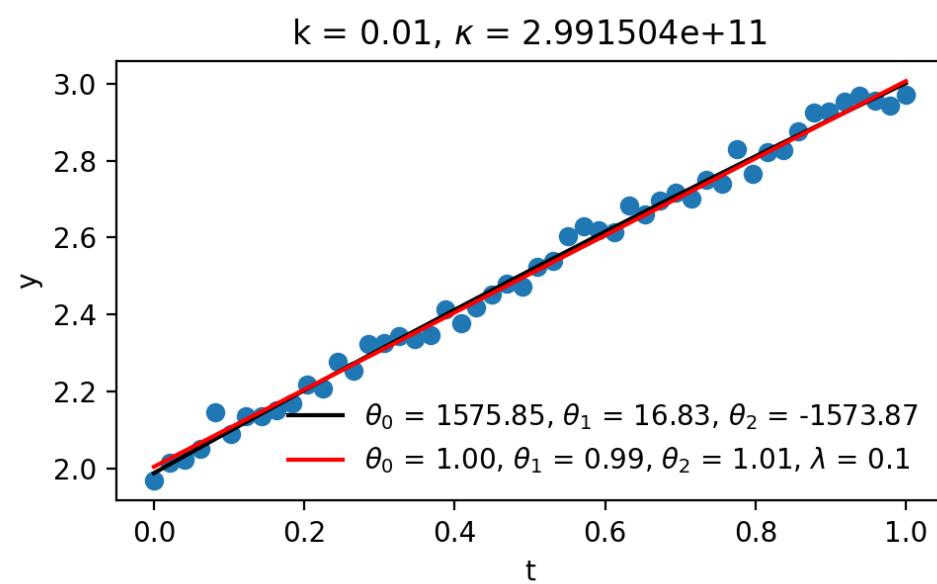
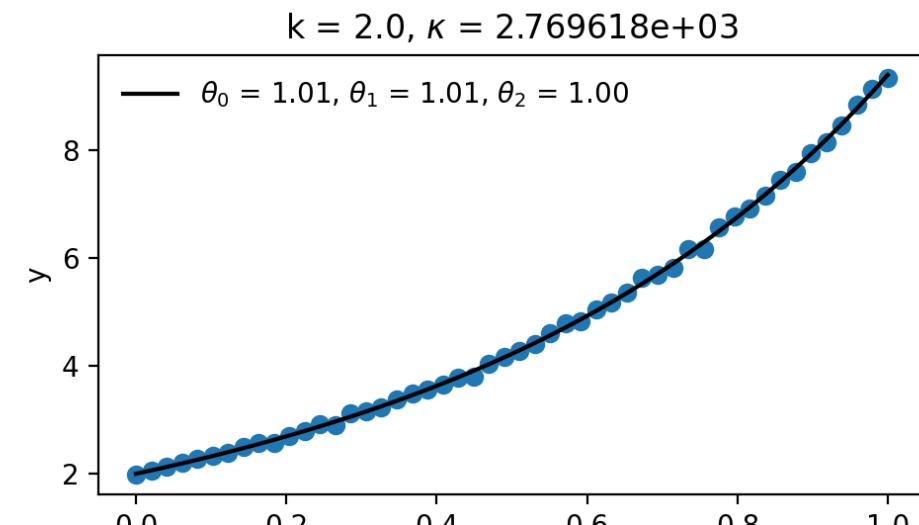
- Another key ML technique is regularisation
- It aims to reduce variance while introducing bias i.e. prevents overfitting

$$\text{MSE} = \text{Var}(f) + \text{Bias}(f)^2 + \sigma^2$$

- Most commonly, Tikhonov regularisation

$$\min_{\theta} \sum_i [y_i - f(x_i, \theta)]^2 + \lambda |\theta|^2$$

- λ is Lagrange multiplier on constraint $|\theta|^2 = 0$
- Other methods include truncated SVD, spectral filtering, parameter priors, and many more



Deconvolution and Regularisation

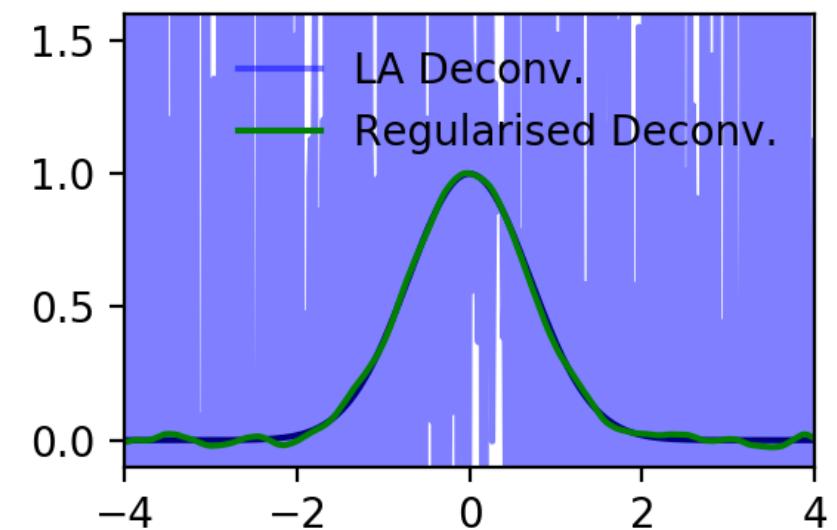
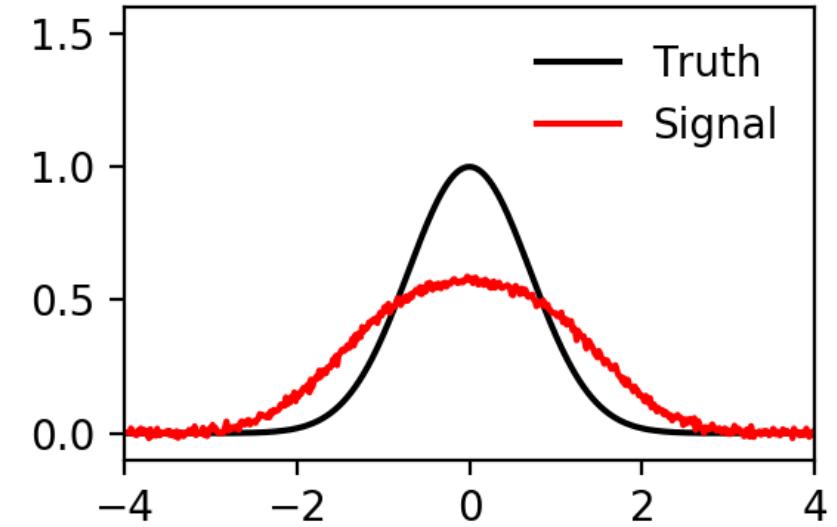
- Convolution is linear

$$F(t) = \int R(\tau - t)f(\tau)d\tau \rightarrow F_i = R_{ij}f_j$$

- However, deconvolution is an ill-posed problem (R_{ij} is *band-limited*)
- Many regularised solution methods for deconvolution
- We will penalise large gradients:

$$\min_{\theta} \sum_i [y_i - R_{ij}f_j]^2 + \lambda (D_{ij}^1 f_j)^2$$

First order finite difference



Non-Linear Least Squares

- Say we have a non-linear model with unknown parameters θ_j :

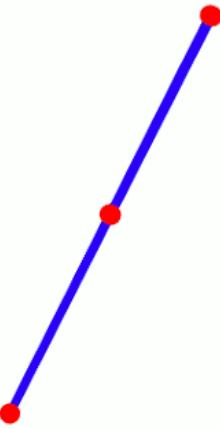
$$f(x, \theta_1 \dots \theta_j \dots \theta_N)$$

- And we are looking for the solution of least squares:

$$\min_{\theta} \sum_i w_i [y_i - f(x_i, \theta)]^2$$

- As you might expect, there is no direct linear algebra solution to this non-linear problem

Double Pendulum at t=0 seconds



Reminder: Newton-Raphson Method

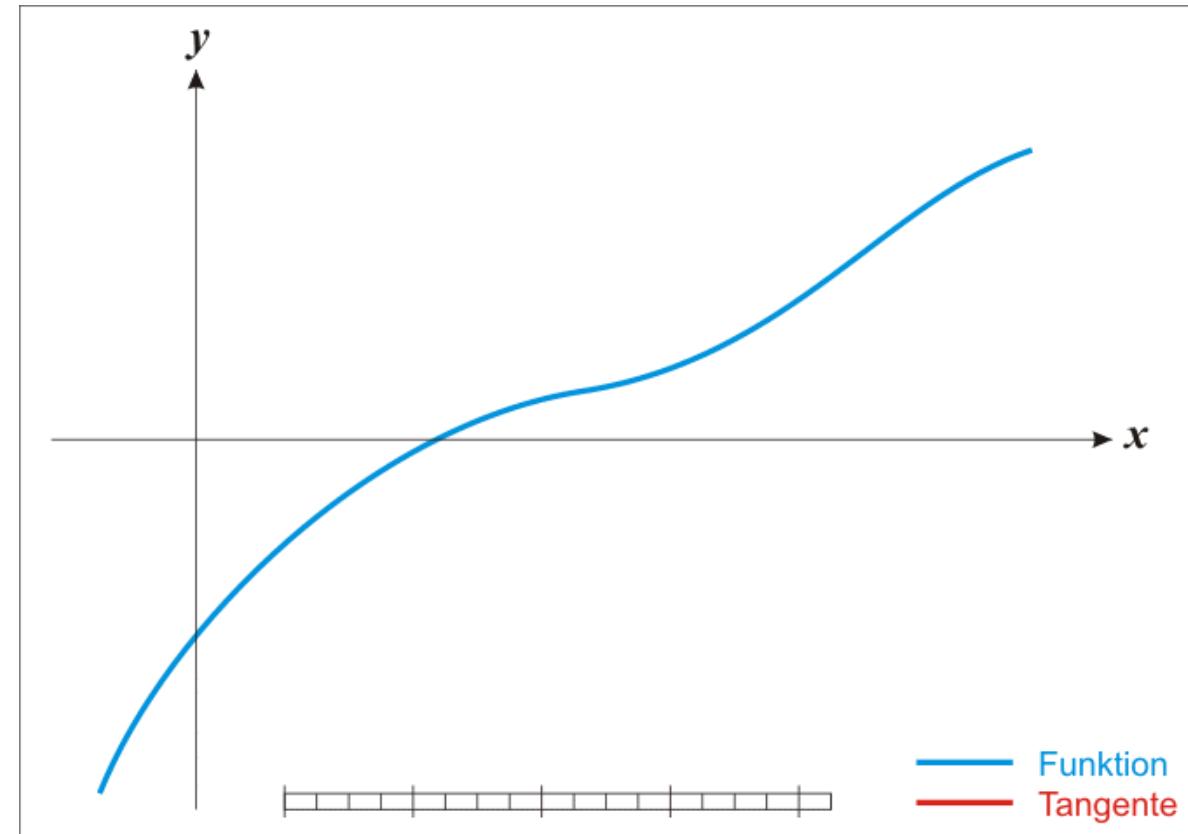


- Iterative method for finding the root of a function:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

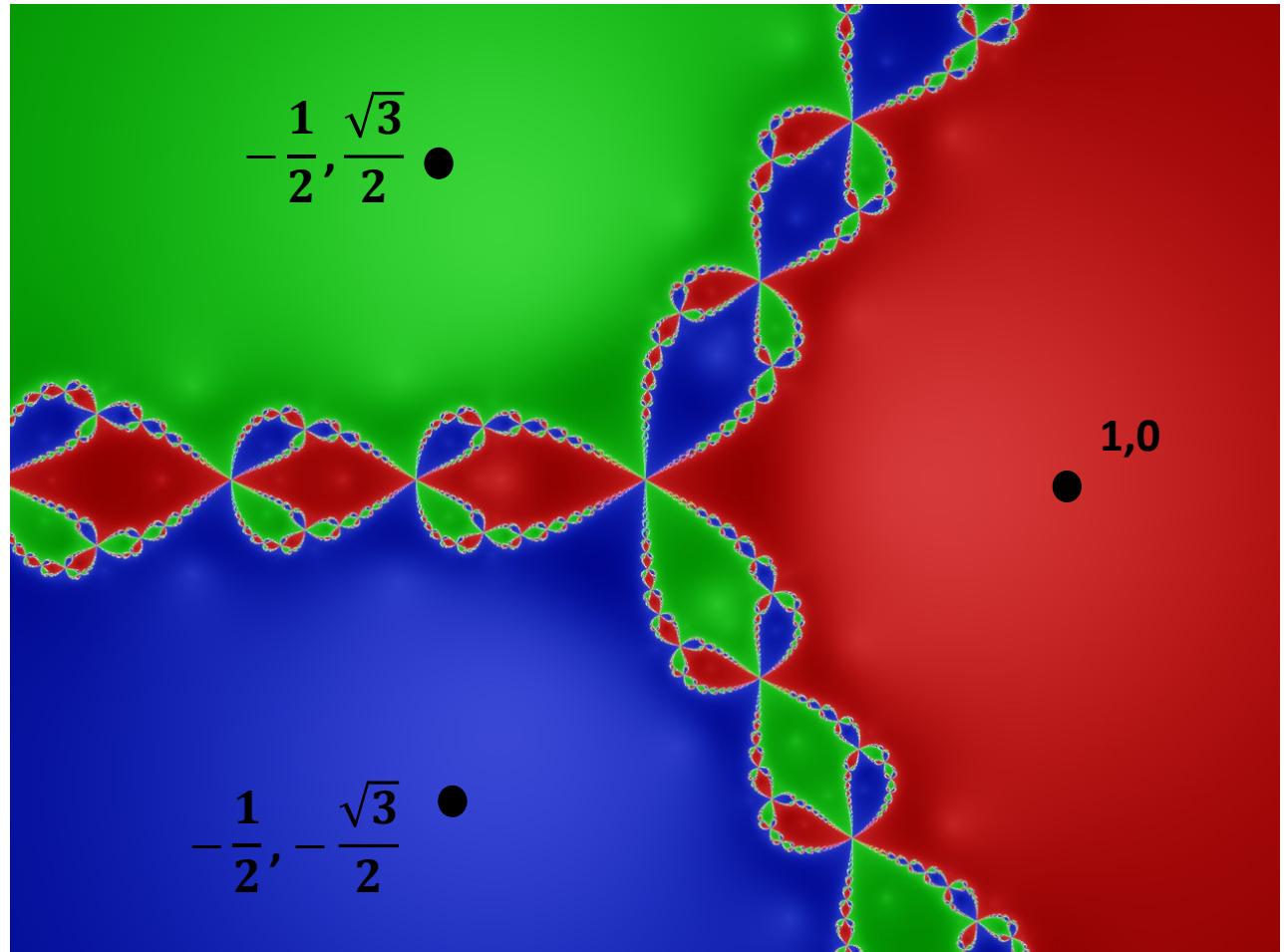
- Minimising a function, $g(x)$, is same as root finding its gradient function, $f(x)$

$$x_{n+1} = x_n - \frac{g'(x_n)}{g''(x_n)}$$



Fun aside: Newton's Fractal

- Depending on starting point, one converges to a different root
- Which root you end up at can have a fractal structure!
- Example: $f(z) = z^3 - 1$
- Example of Holomorphic Dynamics



Non-Linear Least Squares

- We are looking for the solution of least squares:

$$\min_{\theta} S(\boldsymbol{\theta}) = \min_{\theta} \sum_i w_i [y_i - f(x_i, \boldsymbol{\theta})]^2$$

- Lets linearise for small displacements in $\boldsymbol{\theta}$, and drop weights for simplicity:

$$\begin{aligned} & \min_{d\theta} \sum_i [y_i - f(x_i, \boldsymbol{\theta} + d\boldsymbol{\theta})]^2 \\ &= \min_{d\theta} \sum_i [y_i - f(x_i, \boldsymbol{\theta}) + \underline{\mathbf{J}} \cdot d\boldsymbol{\theta}]^2 \end{aligned}$$

- Retrieving a multivariate version of Newton's method:

$$\rightarrow d\boldsymbol{\theta} = (\underline{\mathbf{J}}^T \underline{\mathbf{J}})^{-1} \underline{\mathbf{J}}^T (\mathbf{y} - f(\mathbf{x}, \boldsymbol{\theta}))$$

N.B. we have approximated the Hessian as $2\underline{\mathbf{J}}^T \underline{\mathbf{J}}$, ignoring all 2nd derivatives, so we can get stuck in saddlepoints!

Non-Linear Least Squares

- Retrieving a multivariate version of Newton's method:

$$\rightarrow \underline{d\theta} = (\underline{J^T J})^{-1} \underline{J^T} (\underline{y} - f(\underline{x}, \underline{\theta}))$$

- What if the Jacobian has large condition number? Regularize

$$\rightarrow \underline{d\theta} = (\underline{J^T J} + \lambda \underline{I})^{-1} \underline{J^T} (\underline{y} - f(\underline{x}, \underline{\theta}))$$

- If $\underline{J^T J}$ dominates, then Newton-Raphson step
- If $\lambda \underline{I}$ dominates, then gradient descent step

$$\rightarrow \underline{d\theta} = \frac{1}{\lambda} \cdot \underline{J^T} (\underline{y} - f(\underline{x}, \underline{\theta})) = -\frac{1}{2\lambda} \frac{dS(\underline{\theta})}{d\underline{\theta}}$$

N.B. Levenberg-Marquardt algorithm uses the above but adaptively changes λ , this is default algorithm in `scipy.optimize.curve_fit`

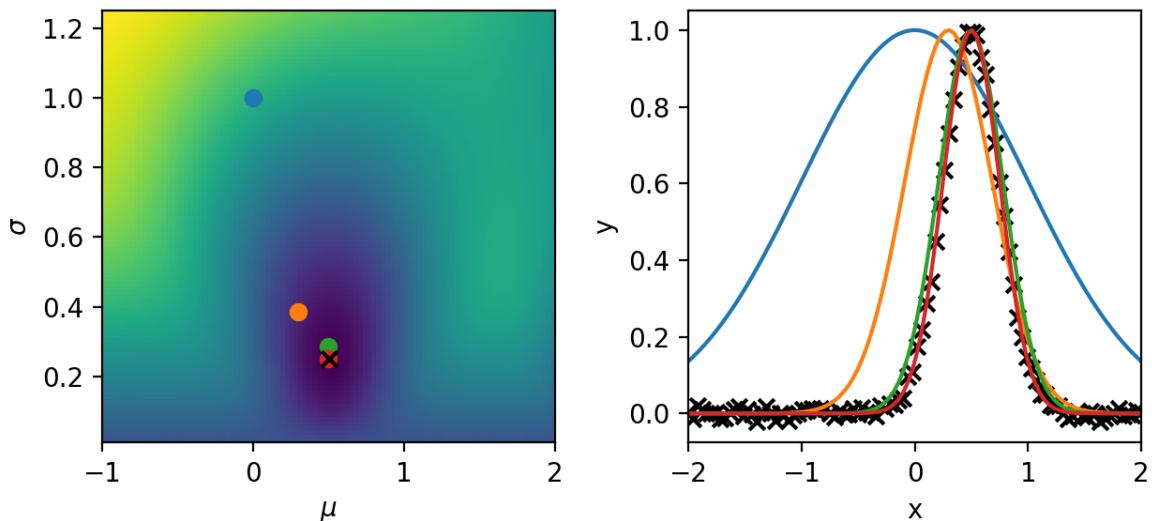
Non-Linear Least Squares

$$\rightarrow d\theta = (\underline{J}^T \underline{J} + \lambda \underline{I})^{-1} \underline{J}^T (y - f(x, \theta))$$

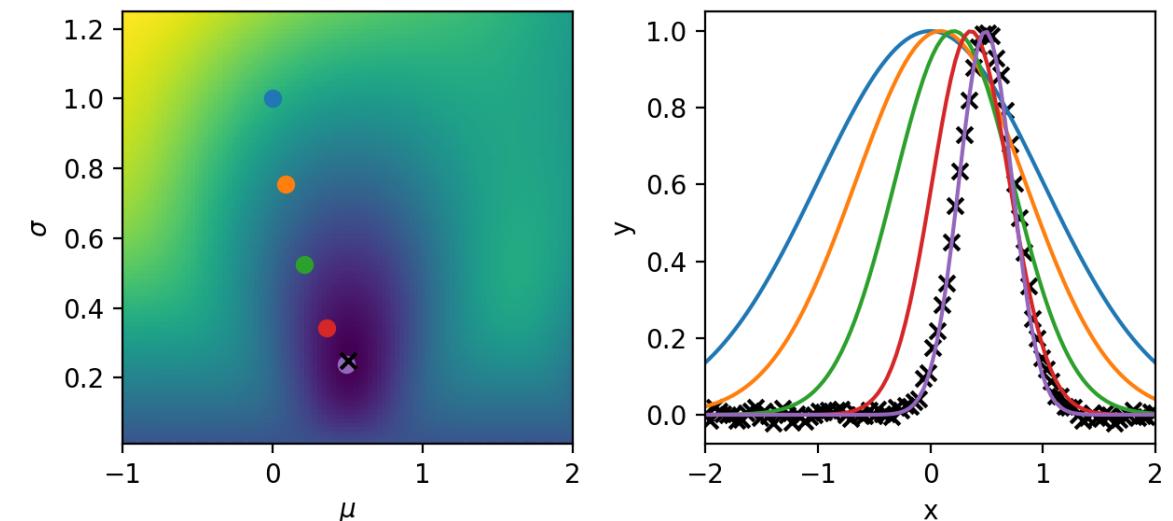
- If $\underline{J}^T \underline{J}$ dominates, then Newton-Raphson step
- If $\lambda \underline{I}$ dominates, then gradient descent step

$$\rightarrow d\theta = \frac{1}{\lambda} \cdot \underline{J}^T (y - f(x, \theta)) = -\frac{1}{2\lambda} \frac{dS(\theta)}{d\theta}$$

Newton-Raphson



Gradient Descent



Differentiable Programming

- One might ask, how do you get the Jacobian, \underline{J} , to perform gradient descent?
- Numerical derivatives: inaccurate and slow for large number of parameters

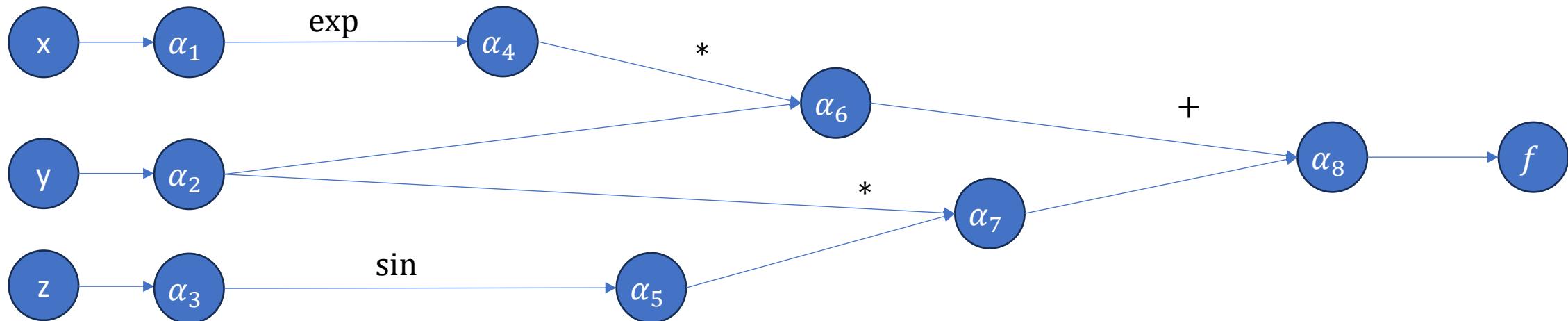
$$\frac{dy}{dx} \approx \frac{y(x + h) - y(x)}{h}$$

- Enter automatic differentiation

Automatic Differentiation Intro

- Every computer program can be written as a computation graph of mathematical ‘primitive’ functions ($+$, x , \exp , \sin , etc.)
- For example:

$$f(x,y,z) = y \exp(x) + y \sin(z)$$



Automatic Differentiation Intro

- In our example:

$$\overline{\alpha_1} = \alpha_4 \overline{\alpha_4} = \alpha_2 \alpha_4$$

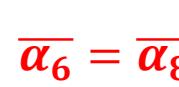
$$\overline{\alpha_4} = \alpha_2 \overline{\alpha_6} = \alpha_2$$



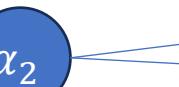
exp



*



$$\overline{\alpha_6} = \overline{\alpha_8} = 1$$



+

Use the adjoint:

$$\bar{x} = \frac{df}{dx}$$



sin



$$\overline{\alpha_5} = \alpha_2 \overline{\alpha_7} = \alpha_2$$

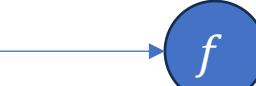
$$\begin{aligned}\overline{\alpha_7} &= \frac{\partial \alpha_8}{\partial \alpha_7} \overline{\alpha_8} \\ &= \overline{\alpha_8} = 1\end{aligned}$$

$$\begin{aligned}\overline{\alpha_3} &= \cos(\alpha_3) \overline{\alpha_5} \\ &= \alpha_2 \cos(\alpha_3)\end{aligned}$$



$$\overline{\alpha_8} = 1$$

+



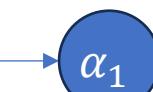
A single 'back-propagation' gives *all* gradients!

Automatic Differentiation Intro

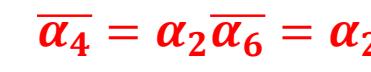
- In our example:

$$\overline{\alpha_1} = \alpha_4 \overline{\alpha_4} = \alpha_2 \alpha_4$$

$$\overline{\alpha_4} = \alpha_2 \overline{\alpha_6} = \alpha_2$$



exp



*

$$\overline{\alpha_6} = \overline{\alpha_8} = 1$$



+



+



f

$$\overline{\alpha_2} = \alpha_4 \overline{\alpha_6} + \alpha_5 \overline{\alpha_7} = \alpha_4 + \alpha_5$$



sin



*



$$\overline{\alpha_8} = 1$$

$$\begin{aligned}\overline{\alpha_3} &= \cos(\alpha_3) \overline{\alpha_5} \\ &= \alpha_2 \cos(\alpha_3)\end{aligned}$$

$$\overline{\alpha_5} = \alpha_2 \overline{\alpha_7} = \alpha_2$$

$$\begin{aligned}\overline{\alpha_7} &= \frac{\partial \alpha_8}{\partial \alpha_7} \overline{\alpha_8} \\ &= \overline{\alpha_8} = 1\end{aligned}$$



A single ‘back-propagation’ gives *all* gradients!

In JAX:

```
1 import jax
2 import jax.numpy as jnp
3 import matplotlib.pyplot as plt
4
5 def f(x,y,z):
6     return y*jnp.exp(x)-y*jnp.sin(z)
7
8 dfdx = jax.grad(f,argnums=(0,))
9 dfdy = jax.grad(f,argnums=(1,))
10 dfdz = jax.grad(f,argnums=(2,))
11
```

Automatic Differentiation Intro

- Gradient descent using Automatic Differentiation (AD) gradients:

$$d\theta = \frac{1}{\lambda} \cdot J^T(y - f(x, \theta)) = -\frac{1}{2\lambda} \frac{dS(\theta)}{d\theta}$$

- Need to:
 1. Define ‘loss function’ $S(\theta)$ which returns a scalar value given input parameters, θ
 2. Using AD framework to compute $\frac{dS(\theta)}{d\theta}$ using AD
 3. Set an initial value of θ
 4. Identify a suitable learning rate – this can be found by checking that $S(\theta)$ is decreasing after θ is updated by $d\theta$

*Other optimisers can use AD gradients,
we shall cover them later*

Using regression for inference

- Finding the ‘best fit’ is an optimisation problem
- In inference, we need additional information – the uncertainties

Bayes Theorem and Likelihoods

- Bayes' Theorem states:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

- Canonical example, medical testing:

$$P(\text{ill}|+) = \frac{P(+|\text{ill})P(\text{ill})}{P(+)} = \frac{\text{Sensitivity} \times \text{Prevalence}}{\text{True Positive} + \text{False Positive}}$$

- Observations always have uncertainty, therefore Bayes rule must be used in parameter estimation



Letter

LII. An essay towards solving a problem in the doctrine of chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a letter to John Canton, A. M. F. R. S

Thomas Bayes

Published: 01 January 1763 | <https://doi.org/10.1098/rstl.1763.0053>

Abstract

Dear Sir, I Now send you an essay which I have found among the papers of our deceased friend Mr. Bayes, and which, in my opinion, has great merit, and well deserves to be preserved.

—

Bayesian Inference

- We wish to find the most likely physical parameters which describe the data:

$$P(\text{model parameters given data}): P(\boldsymbol{\theta}|\mathbf{y}) \propto P(\mathbf{y}|\boldsymbol{\theta})P(\boldsymbol{\theta})$$

- Define a likelihood that *given* the model, $f(\boldsymbol{\theta}, \mathbf{x})$, you would observe the data, say Gaussian errors:

$$P(\boldsymbol{\theta}|\mathbf{y}) \propto P(\boldsymbol{\theta}) \prod_i \exp\left(-\frac{(y_i - f(x_i, \boldsymbol{\theta}))^2}{2\sigma_i^2}\right)$$

- Take the logarithm to separate likelihood and priors

$$\propto \log(P(\mathbf{y}|\boldsymbol{\theta})) + \log(P(\boldsymbol{\theta})) = -\sum_i w_i [y_i - f(x_i, \boldsymbol{\theta})]^2 + R(\boldsymbol{\theta})$$

Maximal posterior = Minimal least squares

Least squares

Regularisation

Pierre-Simon
Laplace



Bayesian Inference – Laplace's Method

- If we expand the log posterior about its optimum:

$$-\log(P(\boldsymbol{\theta}|y)) = F(\boldsymbol{\theta}) \approx F(\boldsymbol{\theta}^*) + \cancel{J_F(\boldsymbol{\theta})} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1}{2} H_F(\boldsymbol{\theta}) (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^2$$

0

A handwritten-style arrow points from the term $J_F(\boldsymbol{\theta})$ to the value 0, indicating it is zero at the optimum.

- We see the posterior is locally a multi-variate normal:

$$P(\boldsymbol{\theta}|y) \sim \exp \left[-\frac{1}{2} H_F(\boldsymbol{\theta}) (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^2 \right]$$

- The Hessian gives us the covariance matrix – this is what scipy's curve_fit is doing

Example in PythonJAX – Differentiable Programming

- Following our non-linear least squares example, lets fit a Gaussian to data with uncertainties estimated using Laplace's method

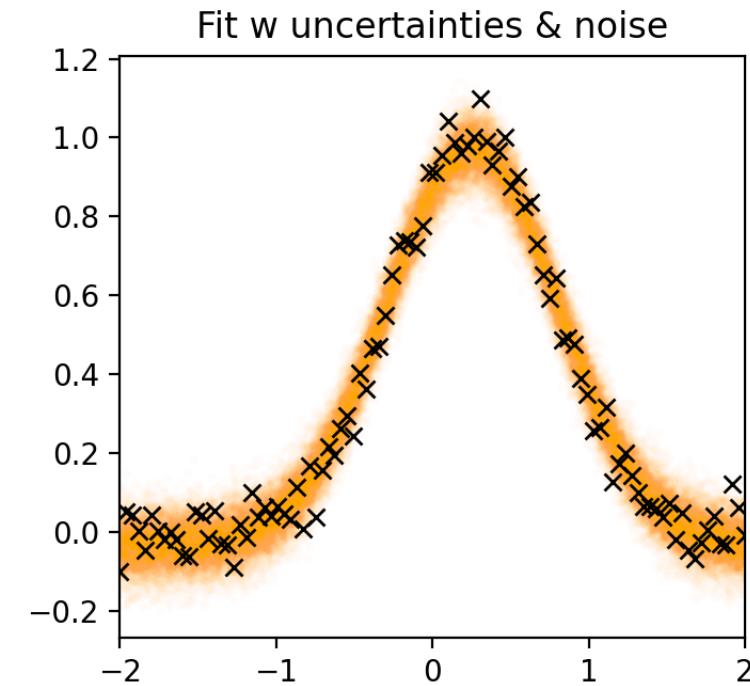
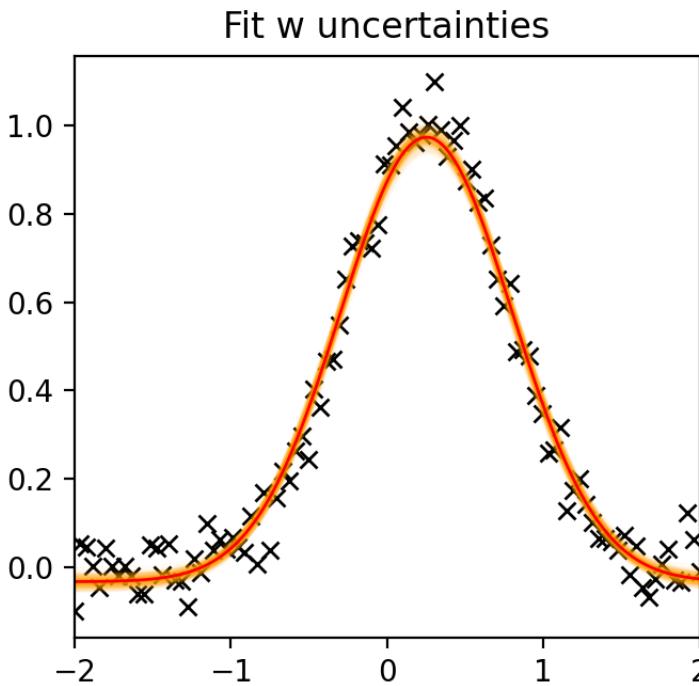
```
def neg_log_likelihood(mu,sig,A,B,x_data,y_data,yerr):
    """ Function to compute the negative log likelihood of our fit to data """
    y_model = Gaussian(x_data,mu,sig,A,B)
    return 0.5*jnp.sum(((y_data-y_model)/yerr)**2)

""" Gradient of neg_log_likelihood function w.r.t. input parameters """
# Here I am doing each gradient explicitly, but it can be done vectorially
# See how Hessian is done
gradL_mu = jax.grad(neg_log_likelihood,argnums=0)
gradL_sig = jax.grad(neg_log_likelihood,argnums=1)
gradL_A = jax.grad(neg_log_likelihood,argnums=2)
gradL_B = jax.grad(neg_log_likelihood,argnums=3)

# N.B. we are going to perform gradient descent so
# the Jacobian of interest is the one with respect to the loss and not the model function
def Jacobian(mu,sig,A,B,x_data,y_data,yerr):
    """ Construct the Jacobian using the Automatic Differential gradient functions """
    J0 = gradL_mu(mu,sig,A,B,x_data,y_data,yerr)
    J1 = gradL_sig(mu,sig,A,B,x_data,y_data,yerr)
    J2 = gradL_A(mu,sig,A,B,x_data,y_data,yerr)
    J3 = gradL_B(mu,sig,A,B,x_data,y_data,yerr)
    J = jnp.array([J0,J1,J2,J3])
    return J

def Hessian(mu,sig,A,B,x_data,y_data,yerr):
    """ Construct the Hessian using the Automatic Differential gradient functions """
    H = jax.hessian(neg_log_likelihood,argnums=(0,1,2,3))(mu,sig,A,B,x_data,y_data,yerr)
    return jnp.array(H).reshape(4,4)

# Truth values of inputs
```



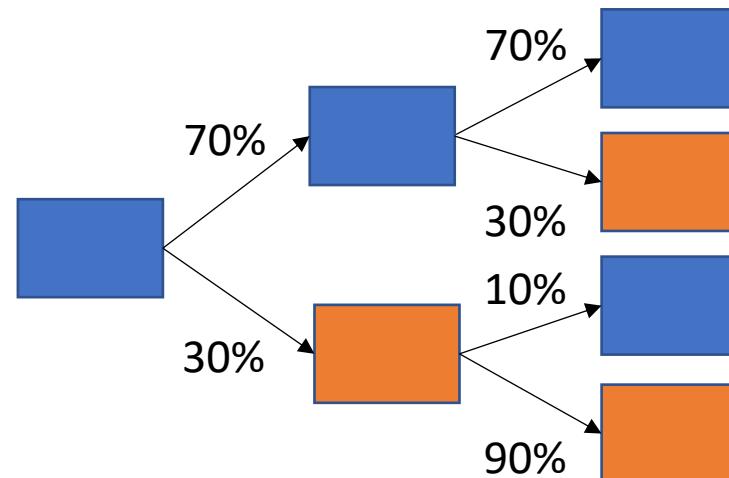
$$\underline{\Sigma} = \underline{H}^{-1}$$

Bayesian Inference: Markov Chain Monte Carlo (MCMC)

- We wish to find the posterior distribution of some parameters but we can only evaluate a function *proportional* to it (Likelihood x Prior)



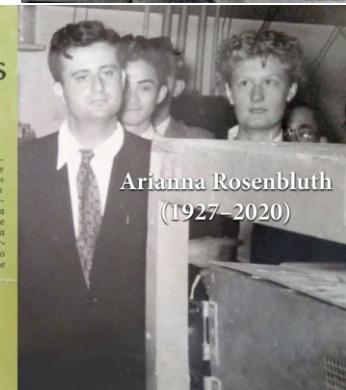
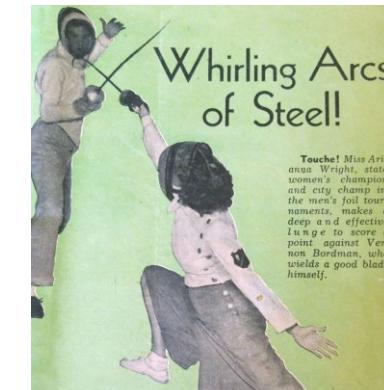
Markov chain:
A sequence of possible events where the next state only depends on the current state



Monte Carlo:
A class of numerical algorithm which involves repeated random sampling

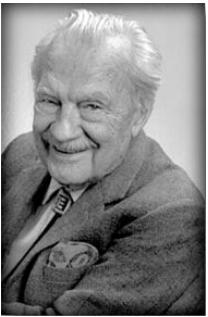


- MCMC: Construct a Markov chain which has an equilibrium distribution of the posterior



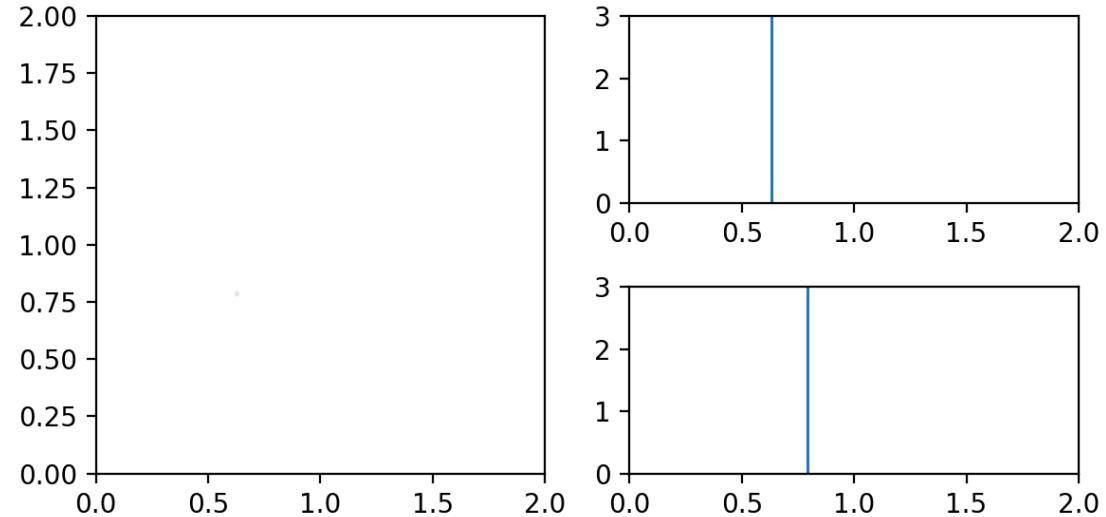
Touche! Miss Arianna Wright, state women's champion and one of the best in the men's foil tournaments, makes a determined effort to touch a point against Vernon Bordman, who wields a good blade himself.

Arianna Rosenbluth
(1927–2020)



Metropolis-Hastings* Algorithm

- $P(x)$ is the probability distribution function we want to sample, $f(x) \propto P(x)$:
 1. Starting at x_i , randomly select a point x^* using jumping function $J(x^*|x_i)$
 2. Compute the acceptance ratio, $\alpha = f(x^*)/f(x_i)$ and a random number $u \in [0,1]$
 - a) If $u \leq \alpha$, accept $x^* \rightarrow x_{i+1}$
 - b) Else, reject and $x_i \rightarrow x_{i+1}$

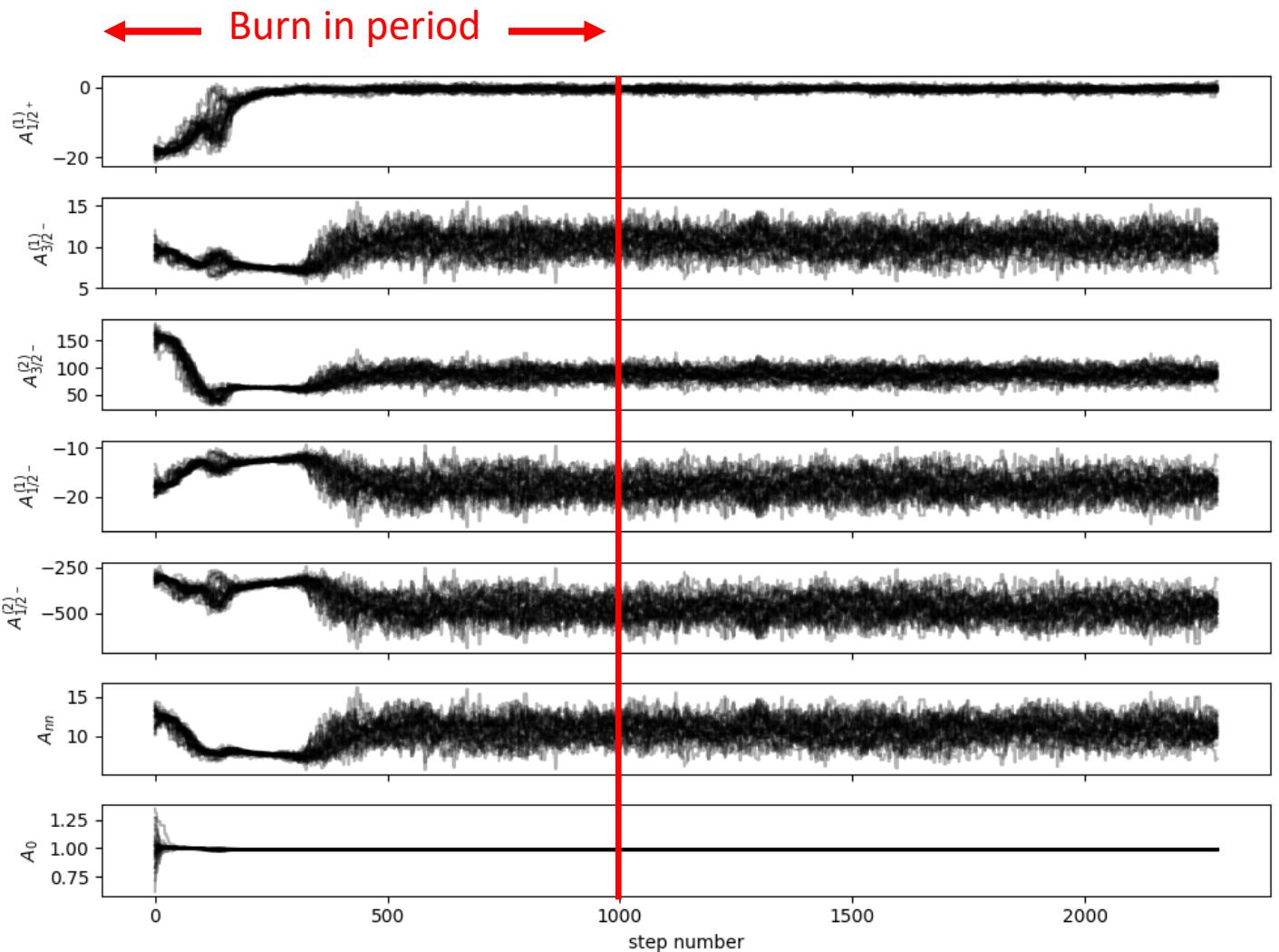


```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def f(x):
5     d = x - mu
6     arg = np.dot(d.T,np.dot(inv_corr,d))
7     return np.exp(-0.5*arg)
8
9 def J(x):
10    y = np.random.multivariate_normal(mean=x,cov=jump_size*np.eye(x.size))
11    return y
12
13 # f parameters, multivariate Gaussian
14 mu = np.array([1.0,1.0])
15 corr = np.array([[0.05,0.02],[0.02,0.05]])
16 inv_corr = np.linalg.inv(corr)
17
18 # Jump function properties
19 jump_size = 0.0025
20
21 # Metropolis algorithm with
22 chain_length = 50000
23 x = 2*np.random.rand(2)
24
25 x1_arr = np.array([])
26 x2_arr = np.array([])
27 u = 1.0
28 alpha = 0.0
29 for i in range(chain_length):
30     while(u > alpha):
31         y = J(x)
32         u = np.random.rand()
33         alpha = f(y)/f(x)
34         x = y
35         u = 1.0
36         alpha = 0.0
37         x1_arr = np.append(x1_arr,x[0])
38         x2_arr = np.append(x2_arr,x[1])
```

* Arguably, the Rosenbluth algorithm

MCMC in practice

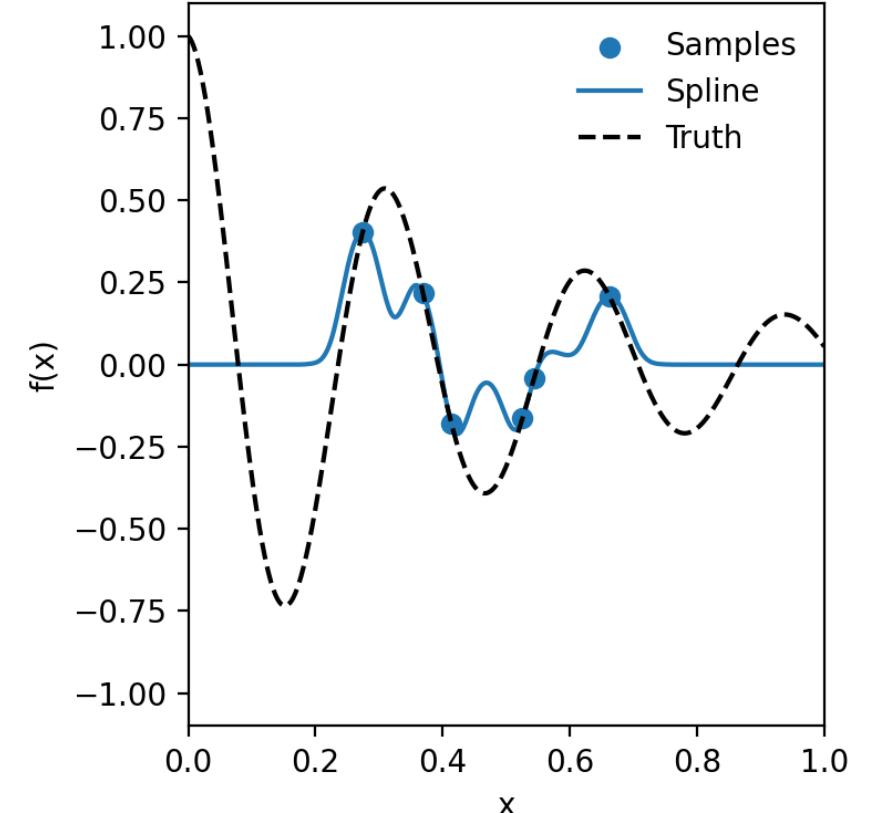
- Good packages: PyMC, *emcee*
- Define your log probability
- Run # walkers for sampling
- Inspect chains and auto-correlation time
- Chop and thin chain for analysis



BREAK

Non-Parametric Regressors

- In the absence of a parametric model, one might want a generic fitting model
- Increasing in various levels of complexity:
 - Nearest-neighbour interpolation
 - Splines
 - Gaussian processes
- Fitting function modifies as you introduce more data? That's a nonparametric regressor



Spline fits are linear least squares:

$$f(x) = \sum_i a_i B(x_i, x)$$

Given data, construct B_{ij} at points x_j
Optimal amplitudes then:

$$\hat{\mathbf{a}} = (\underline{\mathbf{B}}^T \underline{\mathbf{B}})^{-1} \underline{\mathbf{B}}^T \mathbf{y}$$

Conditioning multivariate normals

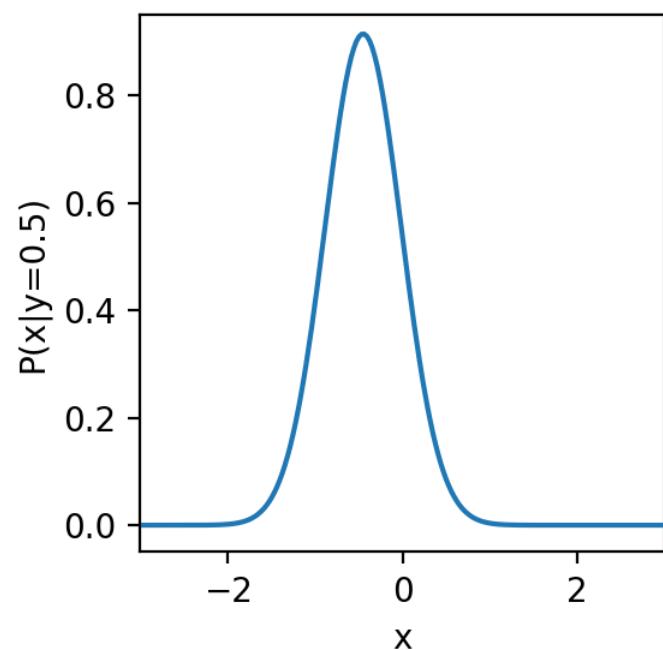
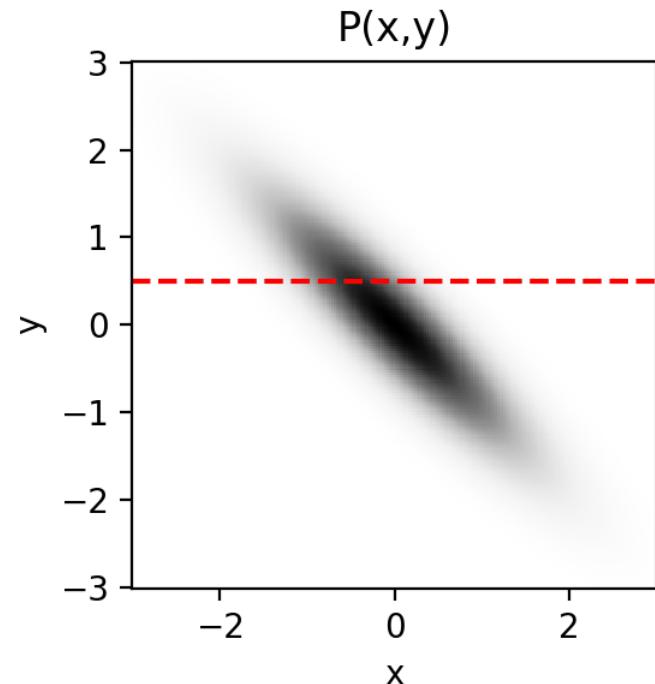
- We wish to predict the random variable Y_0 given data \mathbf{Y}_n , these will be correlated
- Assume joint distribution of prediction + data is multivariate normal:

$$\begin{pmatrix} Y_0 \\ \mathbf{Y}_n \end{pmatrix} \sim N_{1+n} \left[\begin{matrix} \mu_0 \\ \boldsymbol{\mu}_n \end{matrix}, \sigma^2 \begin{pmatrix} R_{11} & \mathbf{R}_{12} \\ \mathbf{R}_{21} & \mathbf{R}_{22} \end{pmatrix} \right]$$

- Conditional distribution of known data has mean and covariance matrix given by:

$$E(Y_0) = \mu_0 + \mathbf{R}_{12}\mathbf{R}_{22}^{-1}(\mathbf{Y}_n - \boldsymbol{\mu}_n)$$

$$\text{Var}(Y_0) = \sigma^2(R_{11} - \mathbf{R}_{12}\mathbf{R}_{22}^{-1}\mathbf{R}_{21})$$



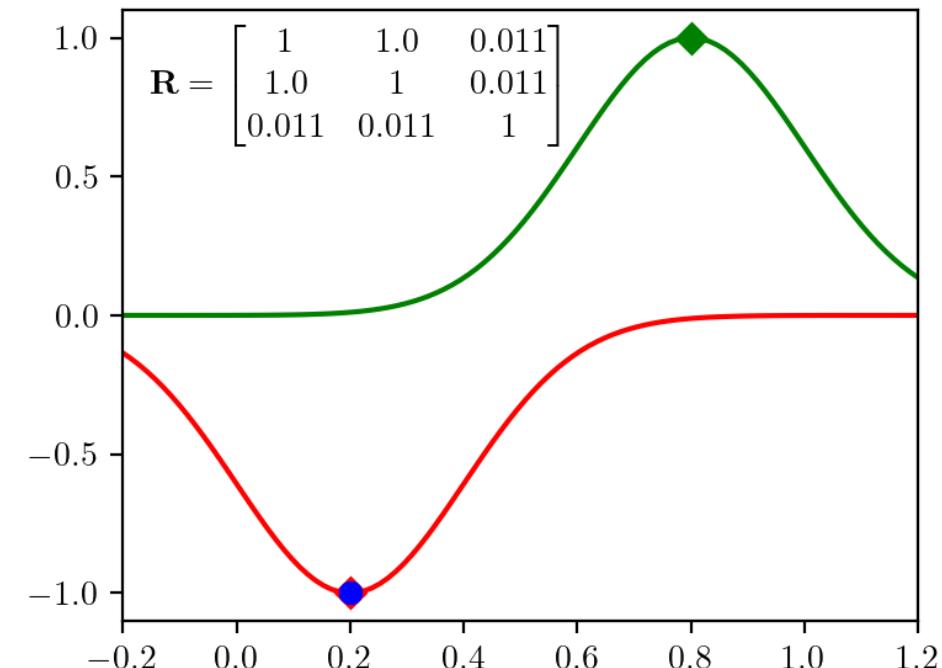
Gaussian Processes

- Gaussian processes use the conditional distribution to predict unobserved points
- Mean function (usually 0) is linear fitting function, $\mu(x) = \sum_i \theta_i f_i(x)$
- Correlation matrix defined by ‘kernel’ function $R(d)$

$$Y(x) = \mu(x) + Z(\mu = 0, R(d))$$

- Correlation function quantifies how correlated points in space are
- How do you pick kernel free parameters...

$$\Sigma = \sigma^2 \begin{bmatrix} R(x_0 - x_0) & R(x_1 - x_0) & \cdots & R(x_N - x_0) \\ R(x_0 - x_1) & R(x_1 - x_1) & \cdots & R(x_N - x_1) \\ \vdots & \vdots & \ddots & \vdots \\ R(x_0 - x_N) & R(x_1 - x_N) & \cdots & R(x_N - x_N) \end{bmatrix}$$



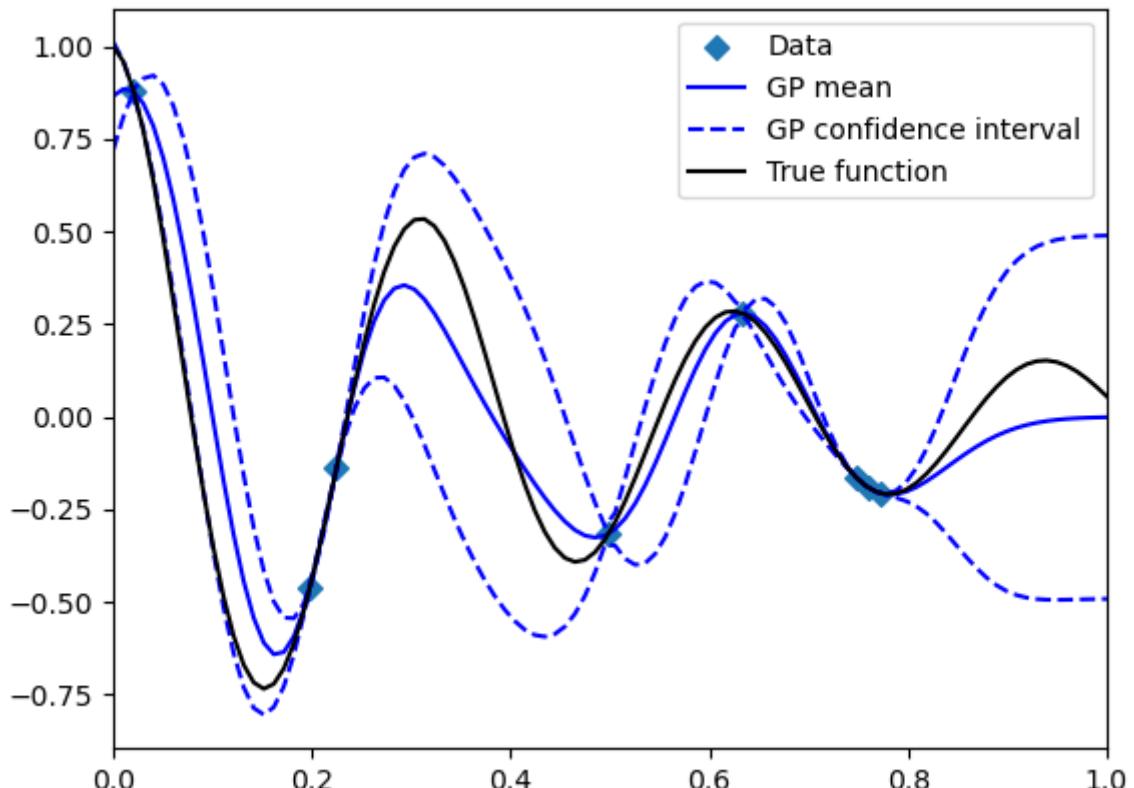
$$E(Y_0) = \mathbf{R}_{12}\mathbf{R}_{22}^{-1}\mathbf{y}_n$$

$$R(d, l) = \exp\left[-\frac{1}{2}\left(\frac{d}{l}\right)^2\right]$$

Gaussian Processes: Example

```
16 def correlation_function(x1,x2,theta):
17     """ Correlation function or kernel of the GP, we will use a Gaussian with length scale theta """
18     return np.exp(-0.5*((x1-x2)/theta)**2)
19
20 def correlation_matrix(x,theta):
21     """ Construct correlation matrix using correlation function """
22     x1,x2 = np.meshgrid(x,x)
23     corr = correlation_function(x1,x2,theta)
24     return corr
25
26 def sigma2_MLE(corr,y):
27     """ Maximum Likelihood estimation of vertical length scale sigma
28         see pg 65-66 of Santner, Williams and Notz"""
29     n = y.shape[0]
30     return np.dot(y,np.matmul(np.linalg.inv(corr),y))/n
31
32 def neg_log_likelihood(theta,x,y):
33     """ Negative Log Likelihood of Gaussian process for a given choice of length scale theta
34         see pg 65-66 of Santner, Williams and Notz"""
35     n = x.shape[0]
36     corr = correlation_matrix(x,theta)
37     sigma2 = sigma2_MLE(corr,y)
38     log_1 = n*np.log(sigma2)+np.log(np.linalg.det(corr))
39     return log_1
40
41 def conditional_distribution(x_trial,x_known,y_known,sigma,theta):
42     """ Calculate p(x_trial | x_known,y_known) for given values of sigma and theta """
43     x_total = np.concatenate((x_trial,x_known),axis=0)
44     N_trial = x_trial.shape[0]
45     corr = sigma**2*correlation_matrix(x_total,theta)
46     corr_11 = corr[:N_trial,:N_trial]
47     corr_12 = corr[:N_trial,N_trial:]
48     corr_21 = corr_12.T
49     corr_22 = corr[N_trial:,:N_trial]
50     # Note unconditional mu = 0
51     corr_22_inv = np.linalg.inv(corr_22)
52     mu_cond = np.dot(np.matmul(corr_12,corr_22_inv),y_known)
53     sig_cond = corr_11 - np.matmul(corr_12,np.matmul(corr_22_inv,corr_21))
54     return mu_cond,sig_cond
55
56 res = minimize(neg_log_likelihood,0.1,args=(X_known,Y_known))
57
58 theta_opt = res.x[0]
59 corr_opt = correlation_matrix(X_known,theta_opt)
60 sigma_opt = np.sqrt(sigma2_MLE(corr_opt,Y_known))
61
62 X_trial = np.linspace(0.0,1.0,100)
63 mu_cond,sig_cond = conditional_distribution(X_trial,X_known,Y_known,sigma_opt,theta_opt)
```

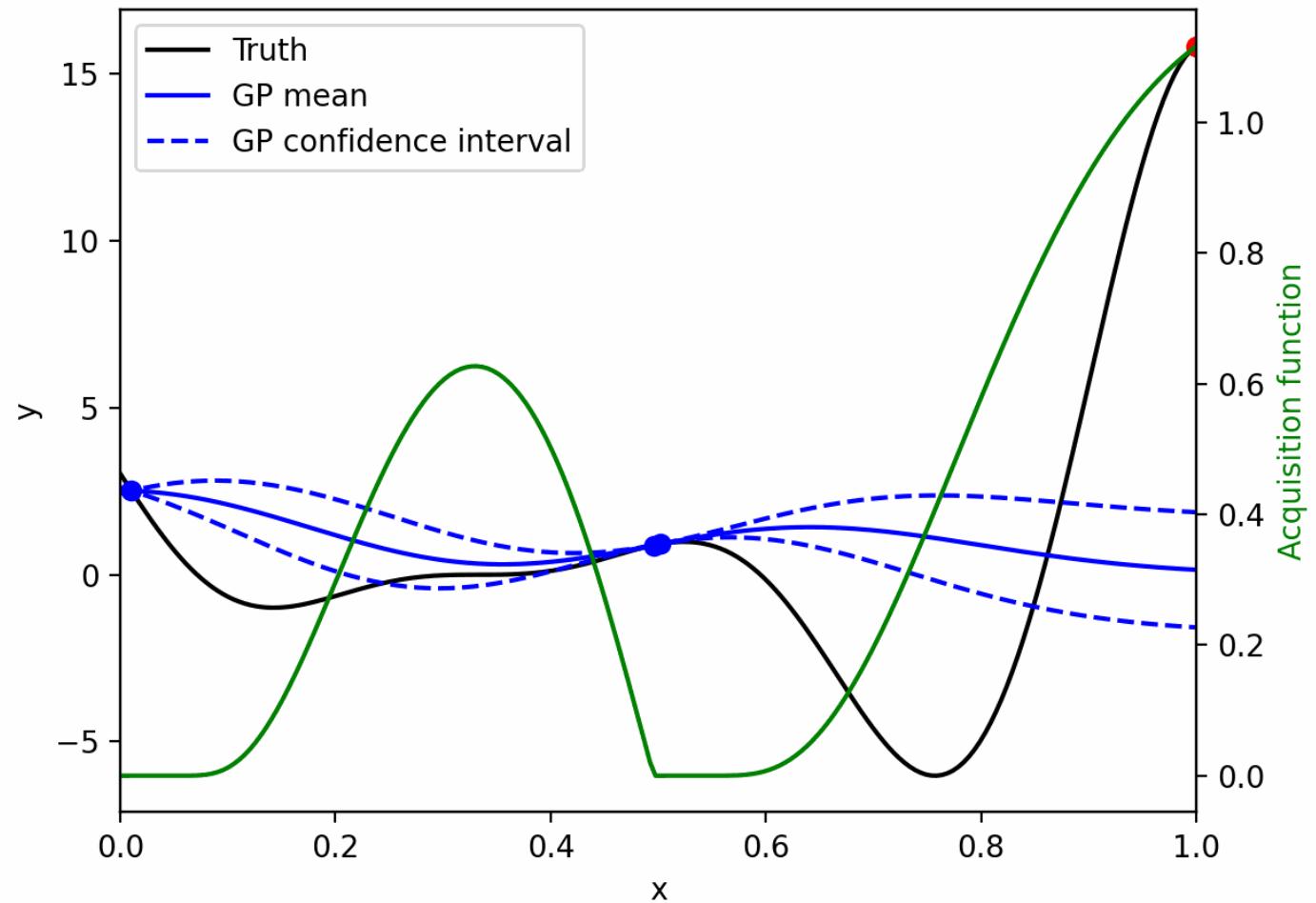
Kernel parameters found by maximal likelihood estimation given known data
Python libraries: *Gpy, sklearn*



Using Gaussian Process: Bayesian Optimisation

One canonical use of Gaussian Processes (GP) is in Bayesian optimisation:

1. From sample points train GP
2. Optimise an ‘acquisition function’ to find next sample point
3. Return to step 1 until global optimum found



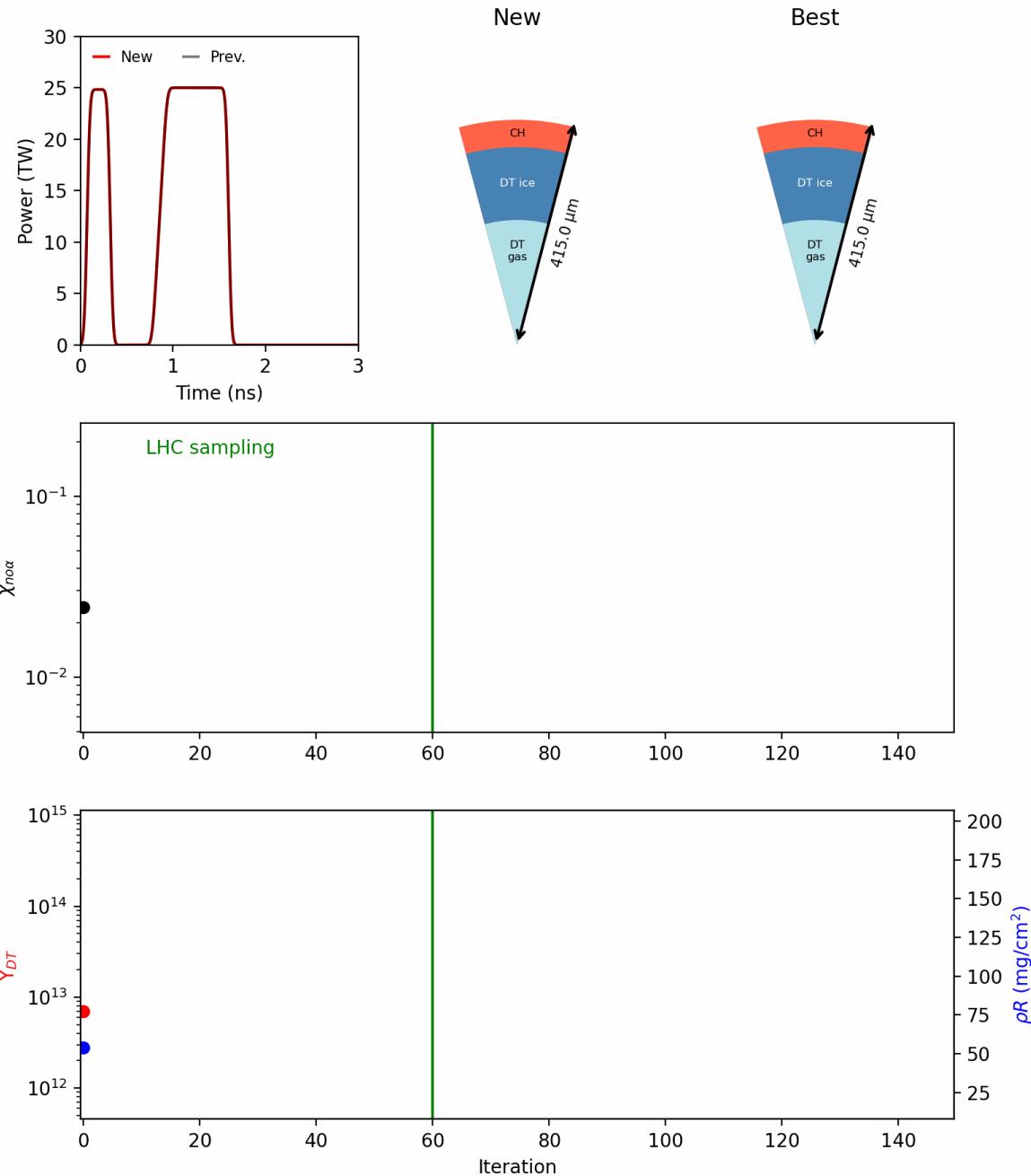
See BayesianOptimisationDemo.py in shared code

Using Gaussian Process: Bayesian Optimisation

In practice:

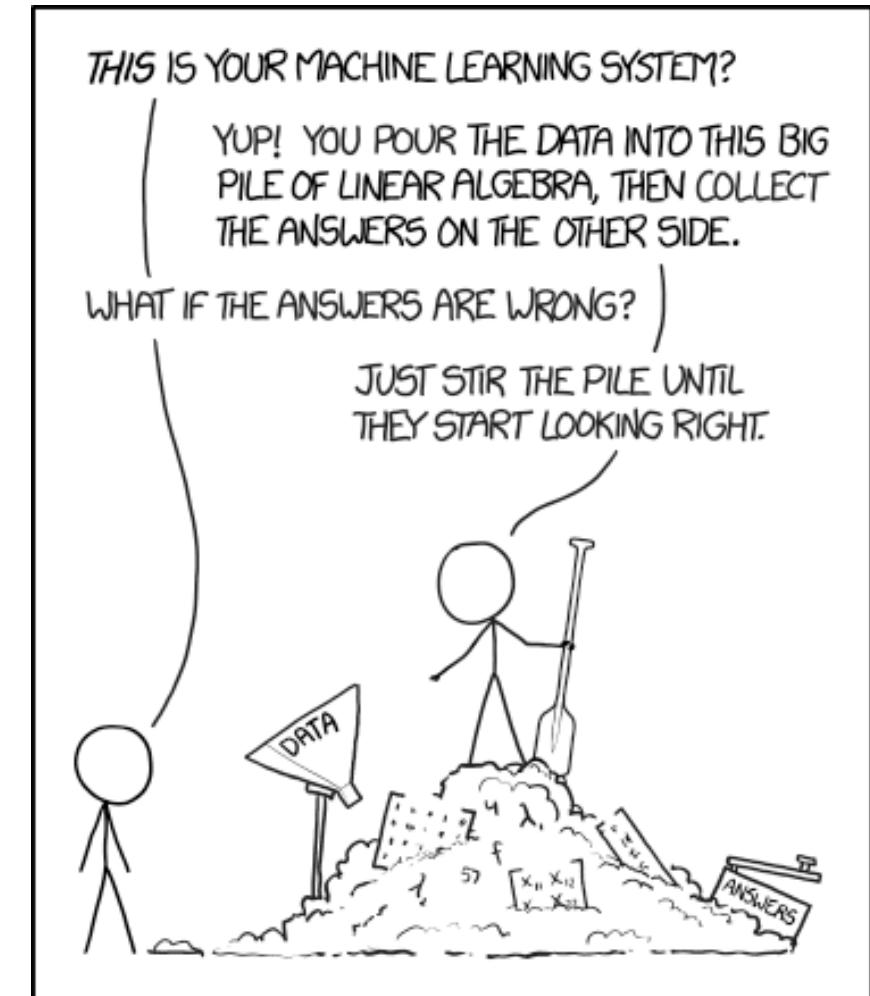
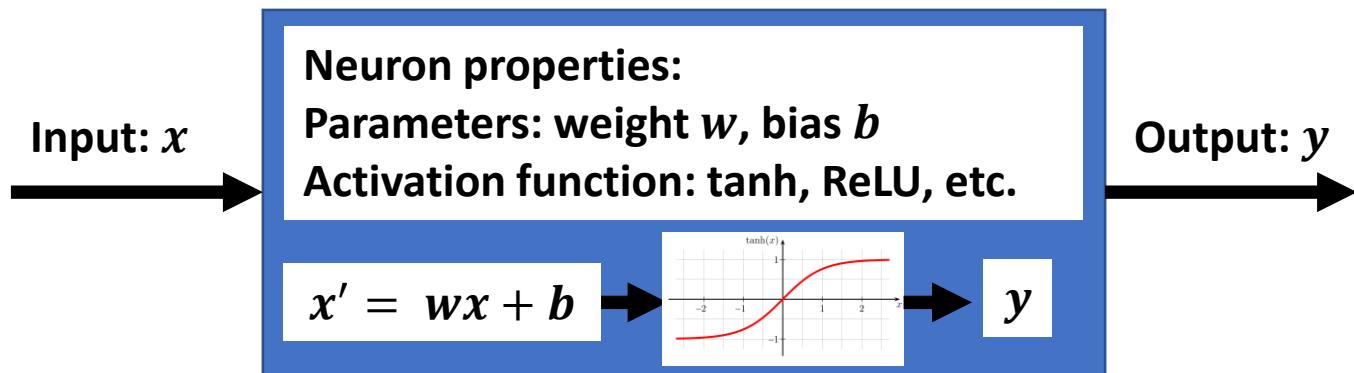
Can be used to optimise
expensive simulations

Usually start with a number of
pseudo-random samples before
training GP and starting BO



Neural Networks

- Neural networks are the canonical AI model
- They are black-box functions with many, many learnable parameters, *ChatGPT 175 billion*
- Many flavours, we will discuss the multi-layer perceptron (MLP)
- Basic unit is the neuron:

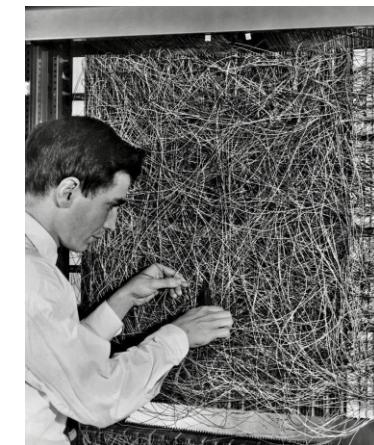
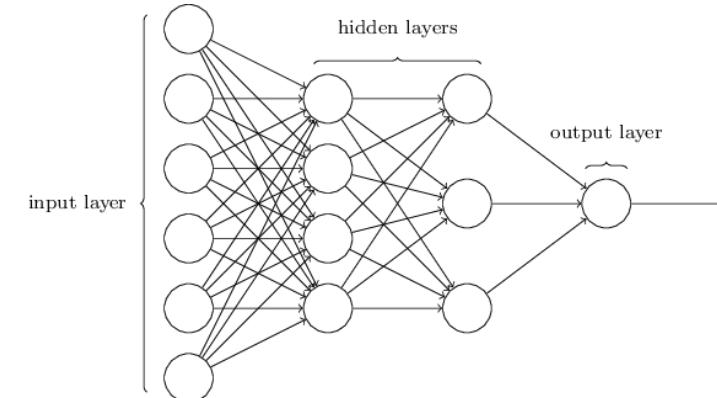


Multi-Layer Perceptron

- Stack neurons into layers and *fully-connect* layers
- Define loss function between predicted and real outputs e.g. mean-squared-error (MSE)
- Update neuron parameters by gradient descent (or variations of it)
- It is just regularised non-linear least squares again!

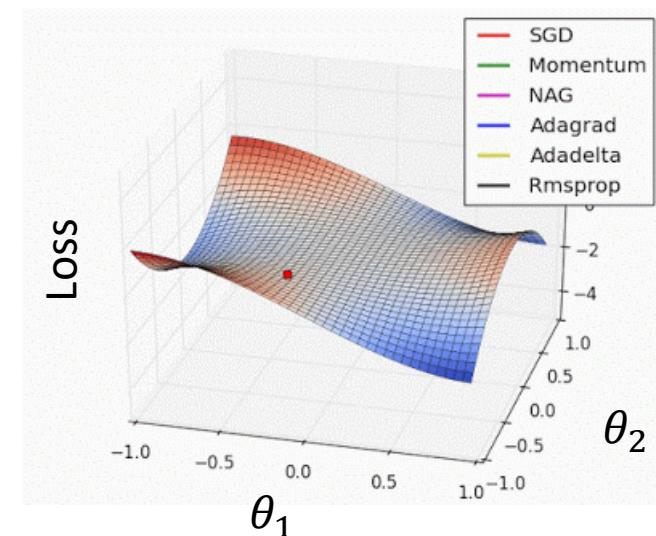
$$d\theta = \alpha \cdot \underline{J^T(y - f(x, \theta))}$$

Learning rate Jacobian found by AD Neural network



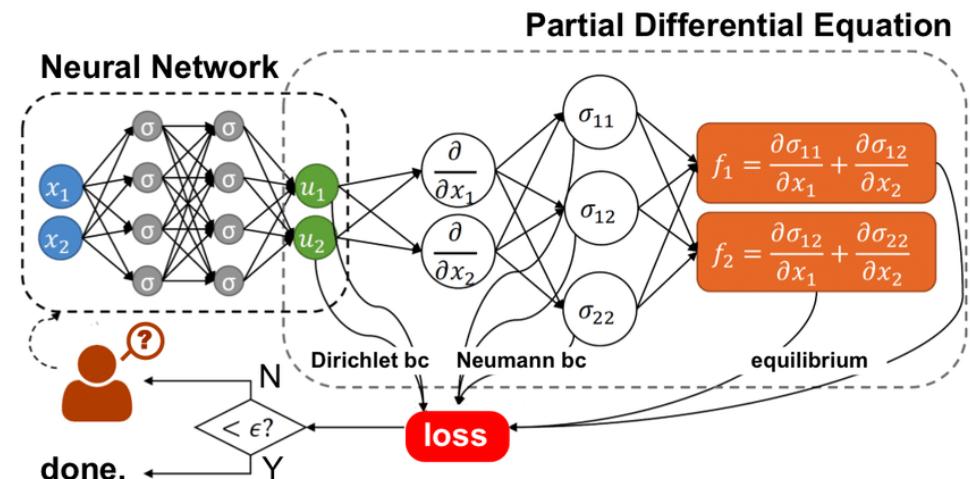
Rosenblatt with
the Mark 1
Perceptron

Various optimisers



Loss terms for neural networks

- The most common are ‘prediction metrics’ such as mean squared error
- Regularisers on weights and biases, common example is ‘L2-norm’ which is the same as Tikhonov, $\lambda|\theta|^2$
 - λ is now hyper-parameter like learning rate, α
- More exotic terms are being introduced, like physics informed neural networks
 - AD allows derivative terms to be added to the loss

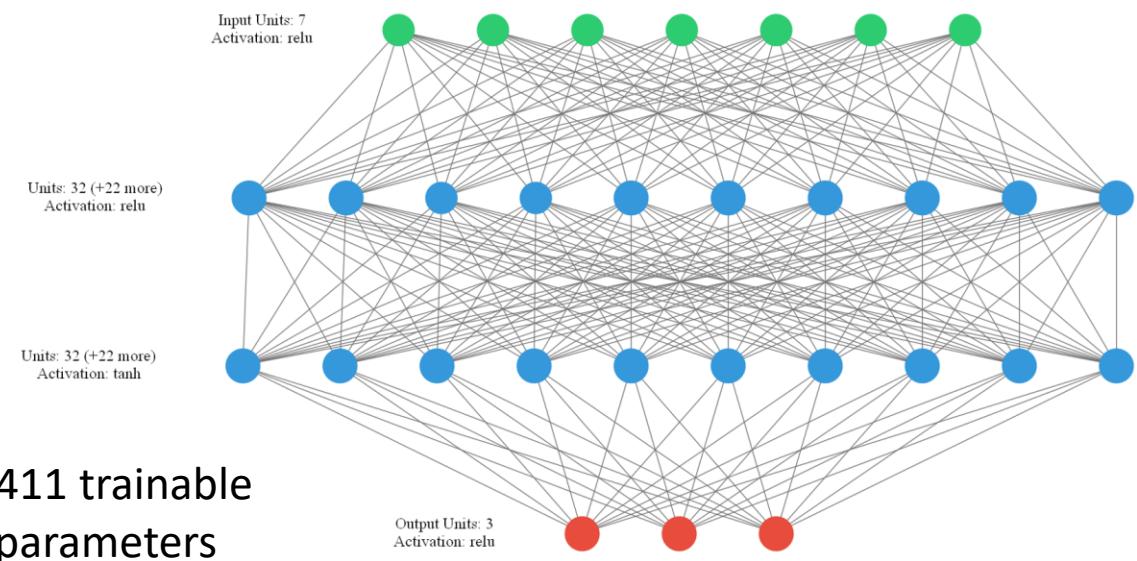
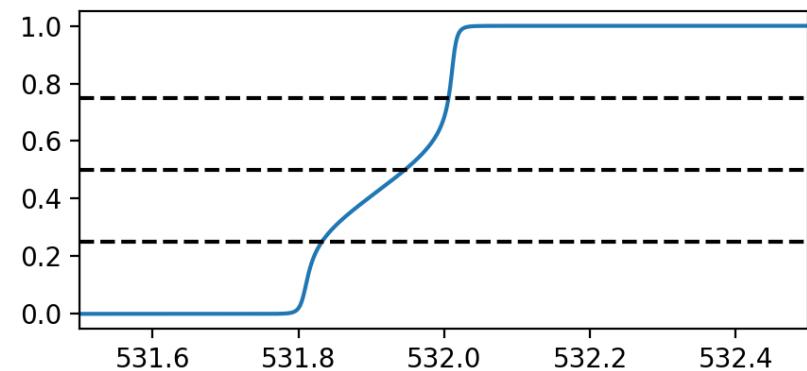
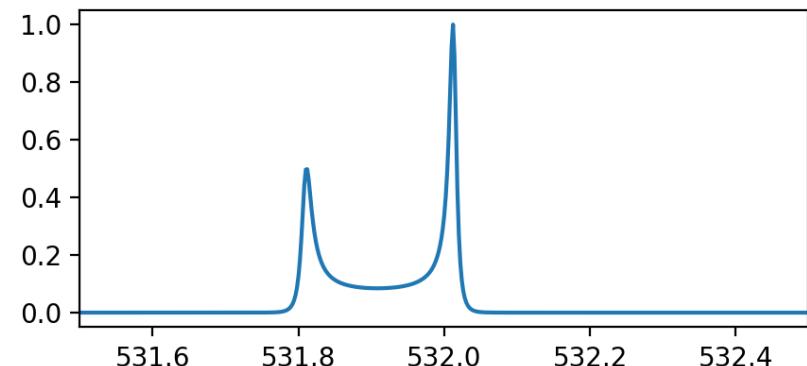


Neural Networks: Example

- Learn relationship between Si IAW Thomson scattering features and physical parameters

$$\theta = (n_e, \mu, T_e, T_i, Z, v_{fi}, v_{fe})$$
$$y = (\Delta\lambda_m, \Delta\lambda_h - \Delta\lambda_l, a_\lambda)$$

- Selection of the NN architecture can be human judgement or algorithmic
 - Tweaking of model parameters becomes tweaking of hyper-parameters
- Data ‘featurisation’ is often key in producing a good model



Neural Networks: Example

- Learn relationship between Si IAW Thomson scattering features and physical parameters

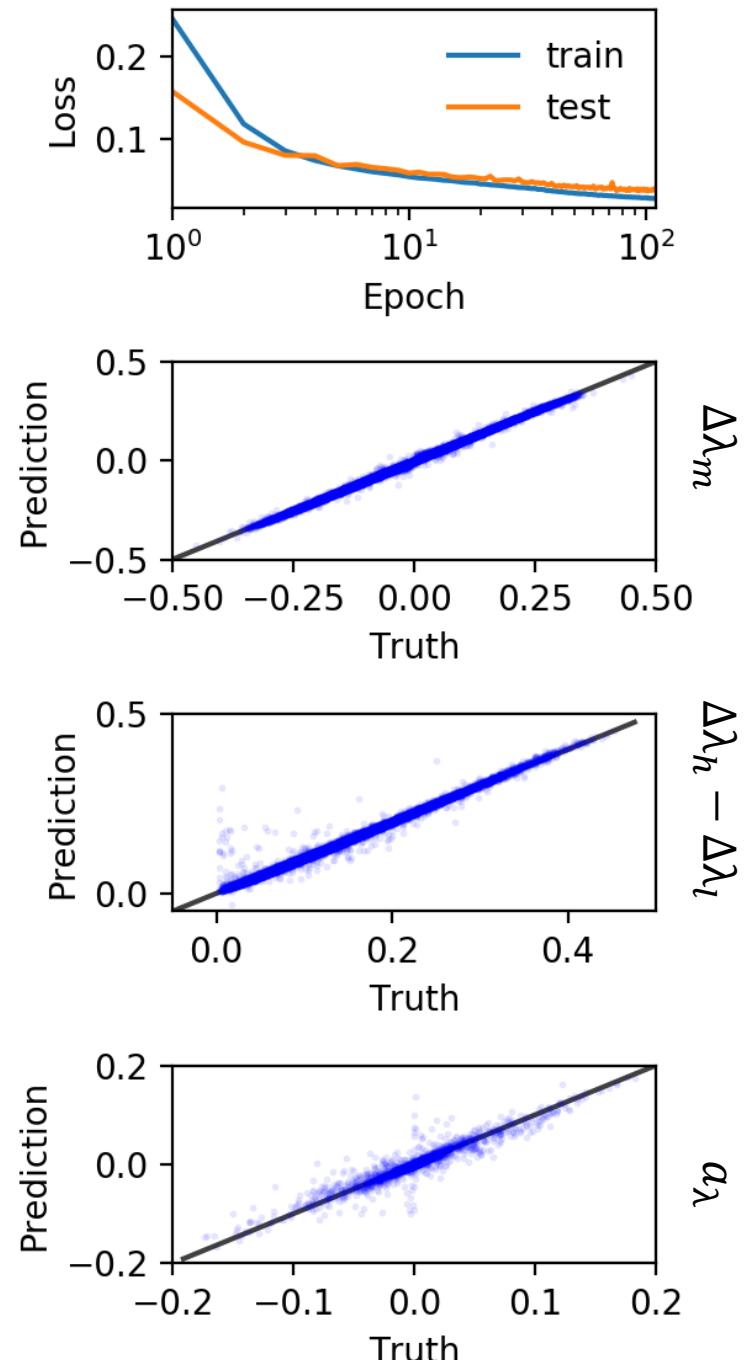
$$\theta = (n_e, \mu, T_e, T_i, Z, v_{fi}, v_{fe})$$
$$y = (\Delta\lambda_m, \Delta\lambda_h - \Delta\lambda_l, a_\lambda)$$

- Deep learning libraries:
keras/tensorflow, sklearn, pytorch

```
# define the keras model
model = Sequential()
model.add(InputLayer(input_shape=(input_size,)))
model.add(Dense(32, input_shape=(input_size,), activation='relu'))
model.add(Dense(32, activation='tanh'))
model.add(Dense(output_size))

model.summary()

model.compile(loss='mean_squared_error', optimizer='adam', metrics=['accuracy'])
history = []
```



Neural Networks: Example

- Learn relationship between Si IAW Thomson scattering features and physical parameters

$$\boldsymbol{\theta} = (n_e, \mu, T_e, T_i, Z, v_{fi}, v_{fe})$$

$$\mathbf{y} = (\Delta\lambda_m, \Delta\lambda_h - \Delta\lambda_l, a_\lambda)$$

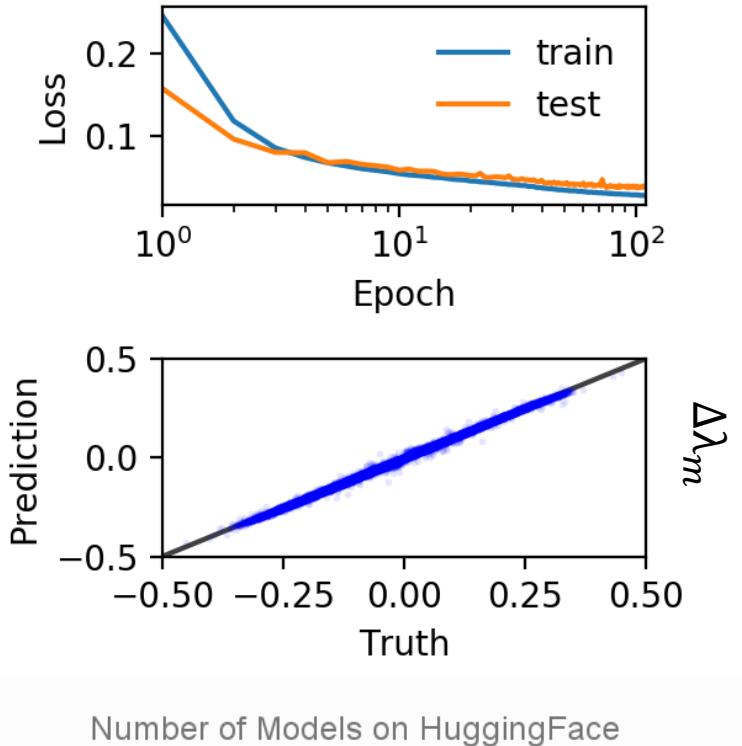
- Deep learning libraries:
keras/tensorflow, sklearn, pytorch

- Learn relationship between Si IAW Thomson scattering features and physical parameters

$$\boldsymbol{\theta} = (n_e, \mu, T_e, T_i, Z, v_{fi}, v_{fe})$$

$$\mathbf{y} = (\Delta\lambda_m, \Delta\lambda_h - \Delta\lambda_l, a_\lambda)$$

- Deep learning libraries:
keras/tensorflow, sklearn, pytorch



Examples in the field

Transfer learning in ICF experiments

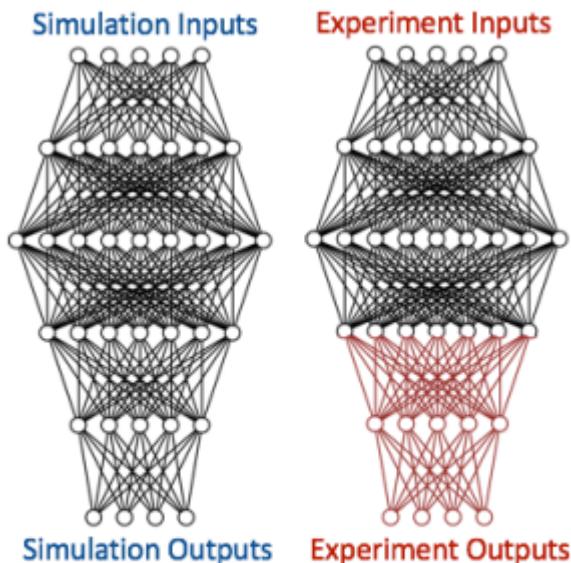
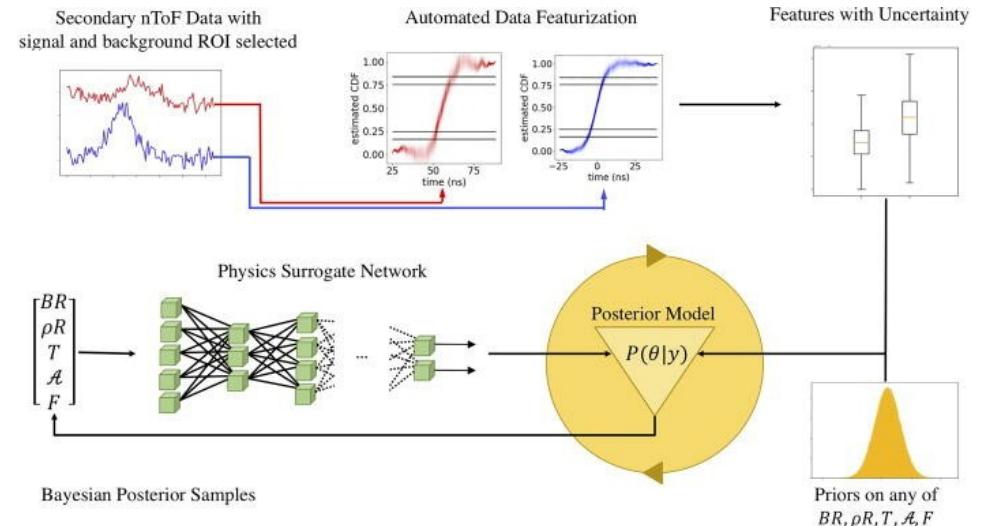
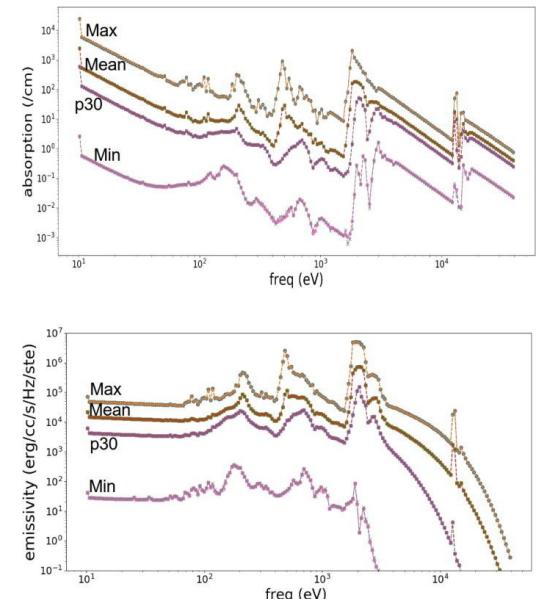
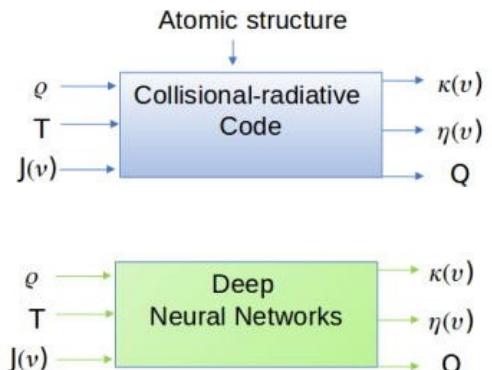


Fig. 1. To transfer learn from simulations to experiments, the first three layers of the simulation-based network are frozen, and the remaining two layers are available for retraining with the experimental data.

Surrogate model of secondaries in MagLIF



NLTE opacity emulators



Method/topic	Resource (Youtube, book, python library, course, etc.)
Machine Learning	<p>https://github.com/josephmisiti/awesome-machine-learning - Curated list of ML resources</p> <p>SciKit-Learn (sklearn), pycaret package gives examples</p> <p>Scipy.optimize</p> <p>“Pattern Recognition and Machine Learning” book</p>
Bayesian Statistics	“Statistical Rethinking” Youtube lecture series and book
Markov Chain Monte Carlo	PyMC, emcee
Gaussian Processes	<p>GPy, GPyTorch, sklearn</p> <p>“Design and Analysis of Computer Experiments” book (WARNING: Mathsy)</p>
Deep Learning	<p>Andrew Ng’s Deep Learning specialisation on Coursera</p> <p>(Can be done for free within time window, otherwise monthly charge)</p>
Neural Networks	Tensorflow, PyTorch, Keras
Differentiable Programming	JAX and associated libraries (diffrax, optax, ...)
Optimisation	Convex Optimisation I and II, Stanford Engineering Everywhere
Software Engineering	<p>Free O’Reilly access: https://www.oreilly.com/library-access/</p> <p>- <i>Python Cookbook</i></p> <p>Imperial’s ‘Essential Software Engineering for Researchers’ course (1 credit)</p>

Technique	Pros	Cons
Traditional optimisation (e.g. least-squares with L-M)	Widely used and easy to start Low number of hyperparameters Familiar techniques	Cannot handle noisy data Often require known Jacobian for best performance Uncertainties use Laplace's method Results depend on starting location Scales poorly Can require regularisation
Gradient free optimisation (e.g. Bayesian optimisation)	Finds global minimum No need for gradient information Can handle noisy data	Scales very poorly Choice of kernel (BO) Spend time exploring domain
Differentiable programming	Gradient information for any program Scales well Allows gradients to appear in loss function	Must write code in differentiable form – harder to debug
Splines, etc.	Flexible functions with local extent Data driven predictor	Often want model parameters not predictor function
Markov Chain Monte Carlo	Explore full posterior distribution Proper Bayesian treatment Does not require gradients (apart from HMC/MC)	Large number of function evaluations needed
Gaussian processes	Like splines but with uncertainties	Scales poorly
Neural networks	Universal function approximators Scales very well Architectures adapted to problems e.g graphs for molecules	Requires large data set Large number of hyperparameters Model training can be expensive

Machine Learning in Plasma Group Research

- Four examples of Plasma Group publications using ML:
 - “Using Differentiable Programming to Learn Closure Relations: An Example in Radiation Transport”
 - “Automation and control of laser wakefield accelerators using Bayesian optimization”
 - “Laser Wakefield Accelerator modelling with Variational Neural Networks”
- These cover Differentiable Programming, Markov Chain Monte Carlo, Gaussian Processes, Neural Networks and optimisation techniques

Learning Closure Relations using Differentiable Programming: An Example in Radiation Transport

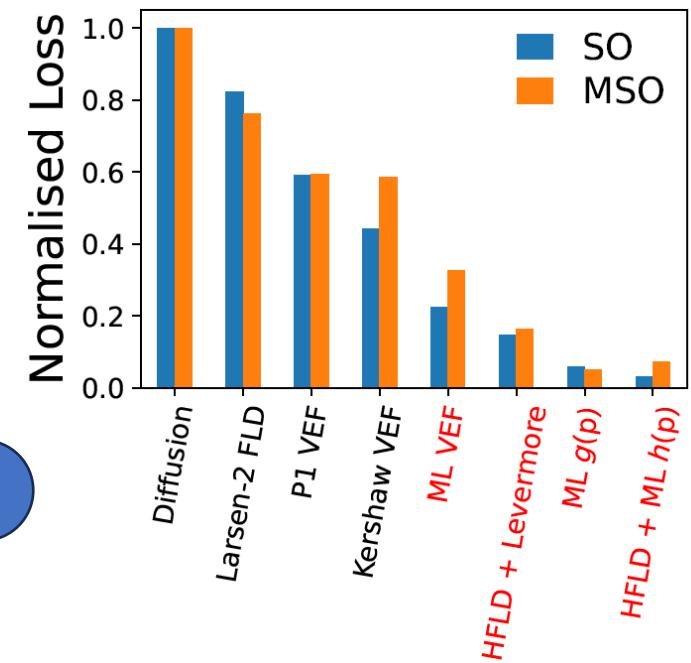
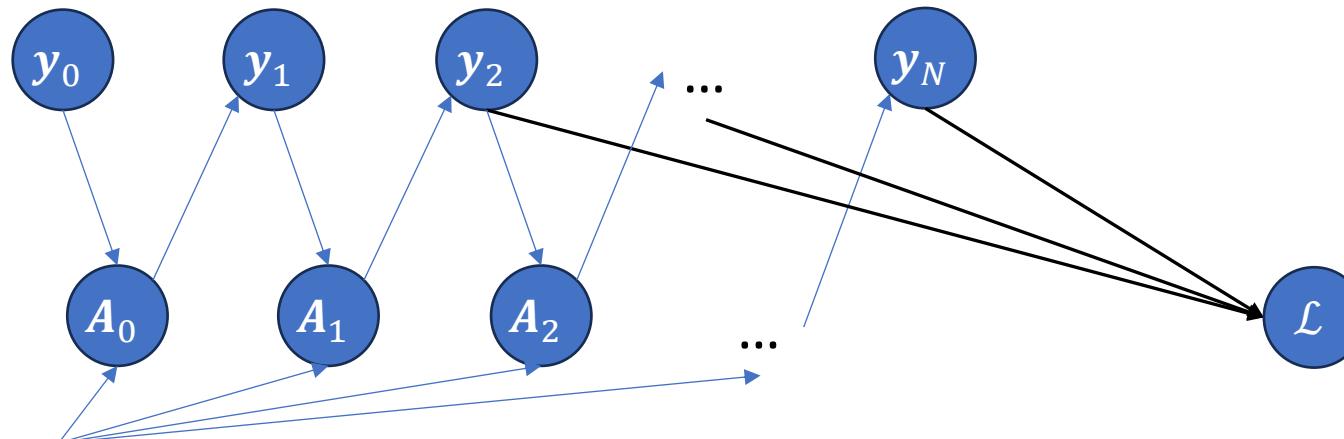
(Crilly et al Accepted at NeurIPS ML 4 Physical Sciences Workshop 2023)

Open-source code



- Can back-propagate gradients through solutions to PDEs
 - *Numerical solutions to PDEs look like pre-trained recurrent graph neural networks*
- Allows closures to be learned from analytic test problems

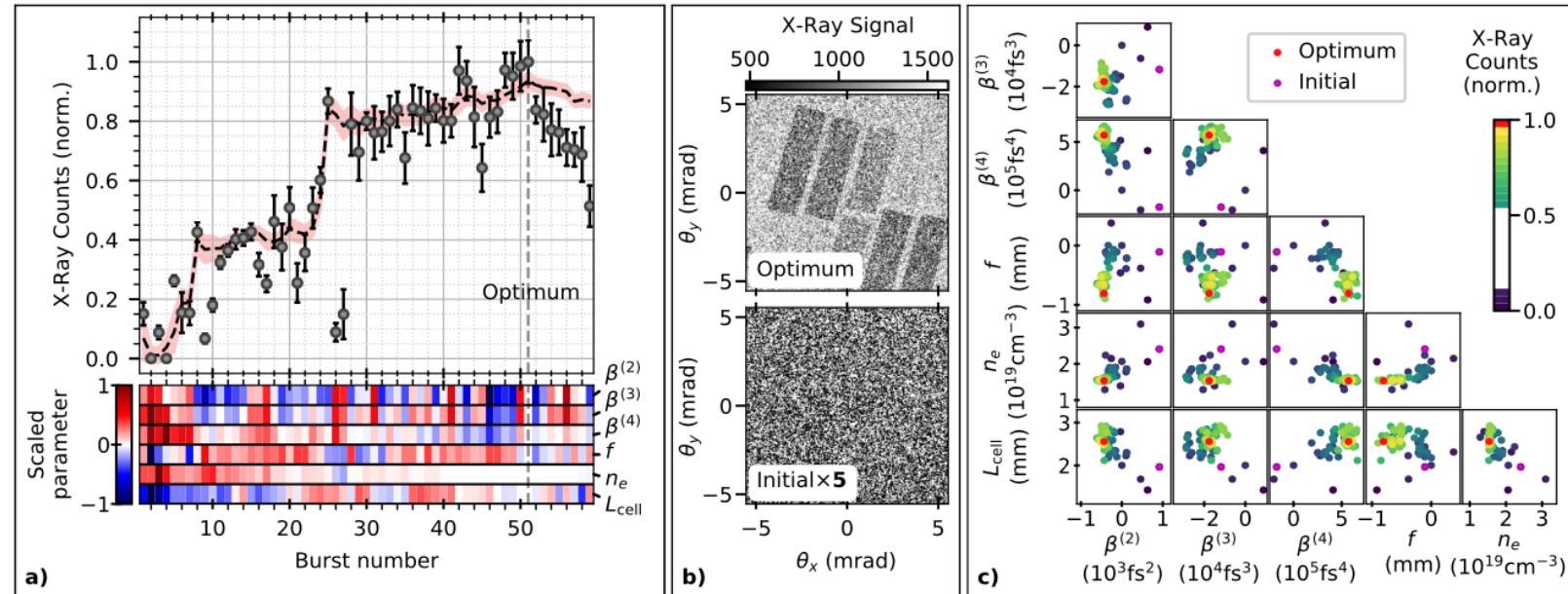
Initial conditions
Model parameters



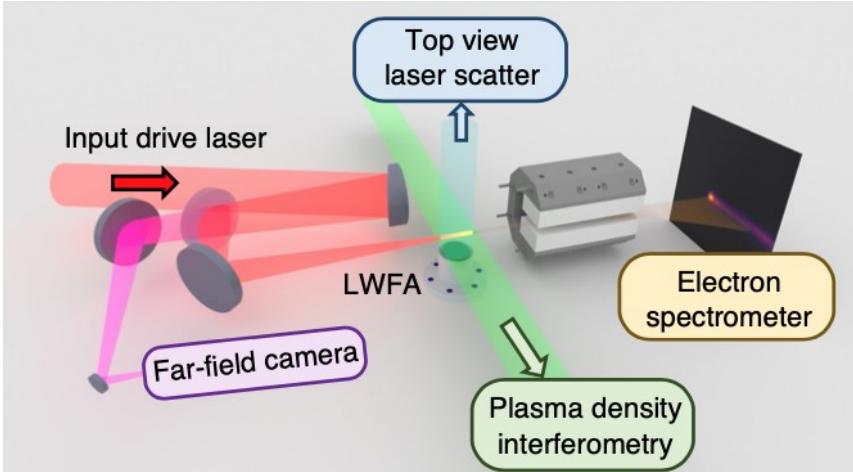
Automation and control of laser wakefield accelerators using Bayesian optimization (Shaloo et al Nat. Comms 2020)

1. Uses Gaussian Process Regression to optimize some property of a laser wakefield accelerator
2. Experimental measurements are made at initial positions
3. GPR model updated with the measurements to form a posterior distribution.
4. An acquisition function is computed and used to select the next measurement location.
5. Steps 3–4 are repeated until the convergence criteria are met.
6. Acquisition function chosen to balance finding the optimum and exploring the function

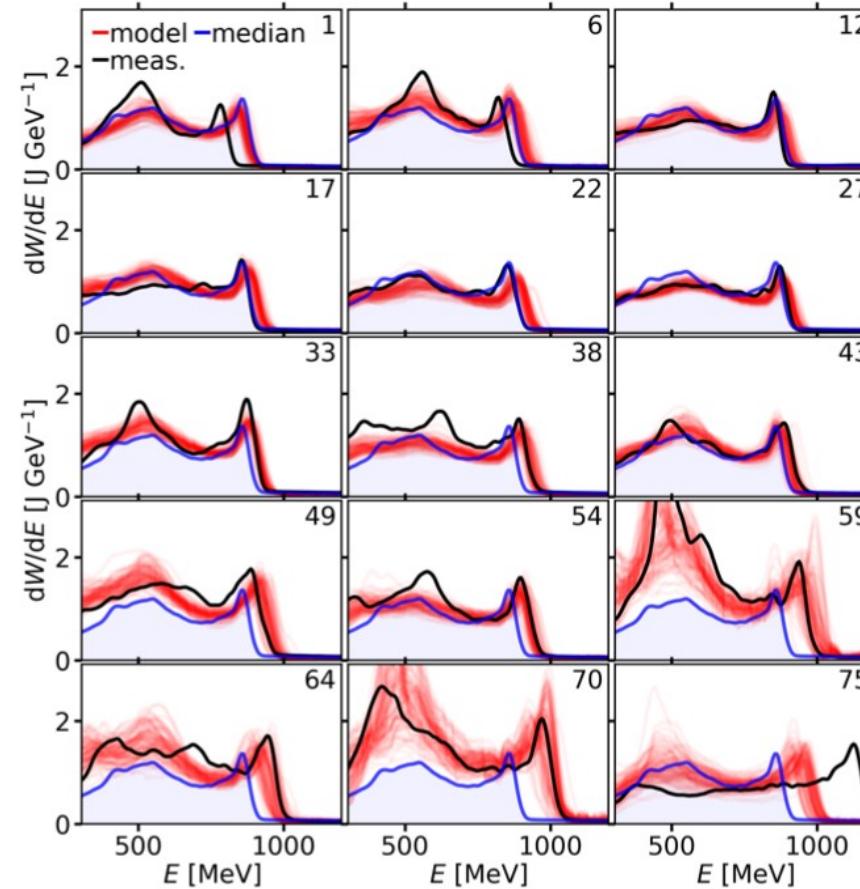
4 laser parameters and 2 plasma parameters controlled to optimize the brightness of generated X-rays



Laser Wakefield Accelerator modelling with Variational Neural Networks (Streeter et al HEDP 2023)



- LWFA spectra differ from shot-to-shot
- In experiments that use the beam we would like to know the spectrum on each shot
 - But if the beam is used, can we use other diagnostics to work out what the beam was?
- We trained a variational neural network based on three non-invasive diagnostics to see if this was feasible
 - Variational neural networks free parameters are random normal variables
 - E.g. Weight $w \rightarrow (\mu_w, \sigma_w)$



15 example predictions – compared to the actual measurement
Black: measured spectrum; Red: model prediction;
Blue: previous best estimate without using a NN

Additional learning oppurtunities

- I am keen to:
 1. Start a book club to focus more on deep learning – Bishop & Bishop
 2. Give a 1-day short course on differentiable simulators (lunch provided)
- Let me know if interested in either or both these ideas!

Links

Please give me some
feedback to improve the
course



Feedback

URL:

<https://forms.gle/37bDPbWVufSjpB9A>



Code and slides

URL:

https://github.com/aidancrilly/ML_Lecture_Demos