



# TIF345 / FYM345

## Advanced Simulation and Machine Learning

### Toward applications and Advanced Regression

Paul Erhart (examiner)

Andreas Ekström

Bernhard Mehlig

Arkady Gonoskov

<https://tif345.materialsmodeling.org>

1

## Course Resources

Course homepage canvas

TIF345 / FYM345 Advanced simulation and machine learning

Below you can find the course-PM for TIF345 / FYM345 "Advanced simulation and machine learning". The main resources for this course, including material for lab assignments as well as additional worked out examples, are maintained in this gitlab repository. To get access to this repository, please create a gitHUB account and send your request to the course examiner ([erhart@chalmers.se](mailto:erhart@chalmers.se)).

Yata (forum)

hypothes.is

### Lecture notes

**tif/fim{345}**  
Advanced Simulation and Machine Learning

1.2. Machine learning

Machine learning (ML) can play an important role in the development of models and analysis of data. Indeed, in the era of abundant and complex data in science as well as industry, from, e.g., particle accelerators, telescopes, and sensors mounted on self-driving cars, the data pattern might only be recognizable through the use of ML.

In science, compared to, e.g., industry, there exists a fundamental difference in regards to the desired foundation of a model. As mentioned above, physicists seek an underlying theory (strategy) to explain and predict future data. ML mostly does the exact opposite. A model derived using ML methods will be agnostic at best, but mostly opaque, when it comes to explaining or understanding of the data. With the ML approach we extract nearly all of the 'intelligence' from data at hand. It remains a challenge for, e.g., supervised learning methods to replace a theory. Indeed, to successfully generalize beyond some training data is limited by the bias-variance effect. Nevertheless, the ML approach is powerful and there exist several cases where the results are excellent. Refs. [2] and [3] provide a more general

Code GitLab

**tif/fim{345}**

For labs jupyterhub

Terminal 2

```
ipython-extensions/tif345/homer: auto -> /root/tif345/homer
```

simple\_machine.ipynb

TIF345 Advanced Simulation and Machine Learning

Lecture 1-2: All models are wrong

General linear models, Bayesian regression, and Gaussian processes

B. Mehlig, ["Machine Learning with Neural Networks: An Introduction for Scientists and Engineers"](#)

# Plan for the next two weeks

- 1-2 Foundations: Linear models, Gaussian processes
- 3-4 Applications: Linear models, Regression, (Gaussian processes)
- 5-7 Neural networks and other ML techniques

L1: Advanced regression

L2: Applications

→ P2a: Regression with cluster expansions

L3: Sensitivity analysis

L4: Global optimization

→ P2b: GPs and Bayesian optimization

# tif/fim{345}

## Advanced Simulation and Machine Learning

### 5. Advanced regression

- 5.1. Ridge regression and beyond
- 5.2. Robust regression
- 5.3. Error correlation
- 5.4. Additional considerations
- 5.5. A brief detour into frequentist statistics
- 5.6. Feature selection and sparse models
- 5.7. Key take-aways

### 6. Toward applications

- 6.1. Alloy cluster expansions
- 6.2. Phonons and force constants
- 6.3. Interatomic potentials

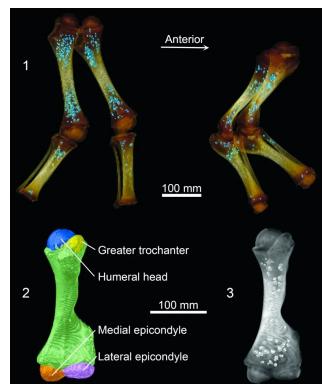
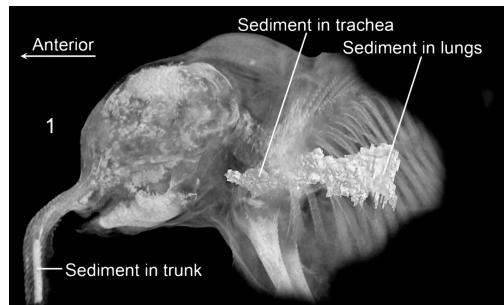
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

Demos 

3

## Applications: Tomography

X-ray tomography

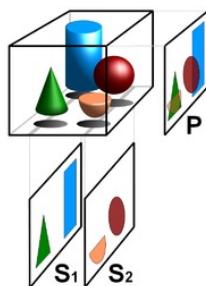


**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

<https://www.livescience.com/16927-baby-mammoth-ct-scans.html>

# Applications: Tomography

**Working principle**  
generate 3D image  
from series of 2D data



**Radon transform (2D)**

$$\mathcal{R}[f(x, y)](H, \phi) =$$

$$\int_{-\infty}^{\infty} f(H(\sin \phi, \cos \phi) + s(\cos \phi, \sin \phi)) ds$$

Function on x-y plane, e.g., density

→ Can be cast as **generalized linear model**



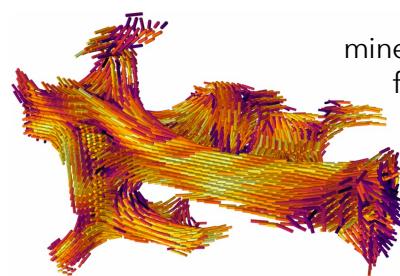
What else could  $f$  be?

Pictures from Wikipedia,  
<https://imgflip.com/gif/29e1g4>

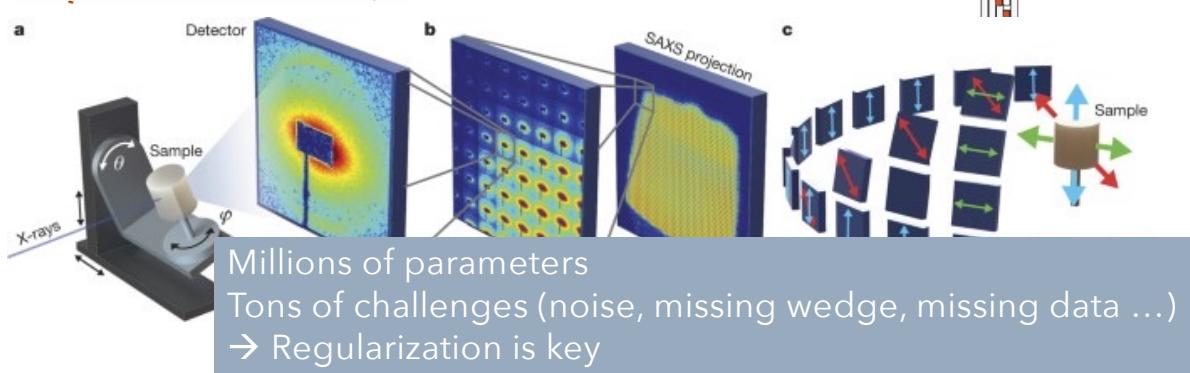
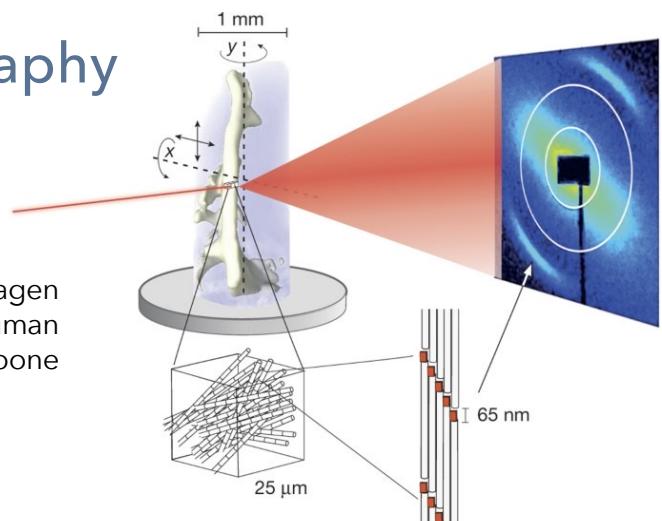
5

# Applications: Tomography

From scalar to tensor tomography  
→ Orientation distribution



mineralized collagen fibrils in a human trabecula bone



# Applications: Atomic scale models

Task: Map positions to potential energy surface (or other properties)

- Interatomic potentials  $\{\mathbf{R}_i\} \rightarrow E \quad \mathbf{f}_i = -dE/d\mathbf{R}_i$

Generalized linear models

- Force constant expansions  $\{\mathbf{u}_i\} \rightarrow E \quad \mathbf{u}_i = \mathbf{R}_i - \mathbf{R}_i^{(0)}$   
 $\mathbf{f}_i = -\Phi_{ij}\mathbf{u}_j - \Phi_{ijk}\mathbf{u}_j\mathbf{u}_k \dots$   
↑ Displacements
- Cluster expansions  $\{\sigma_i\} \rightarrow E$   
↑ Occupation vector

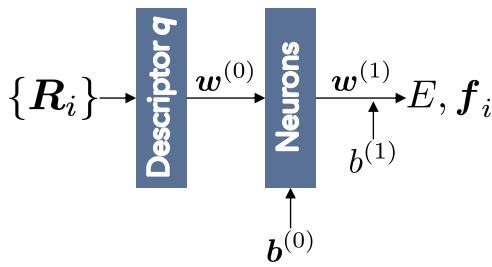
Hundreds to tens of thousands parameters  
 → Regularization is key

**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

7

# Applications: Atomic scale models

Example: Neuroevolution potential (NEP) models



Loss function

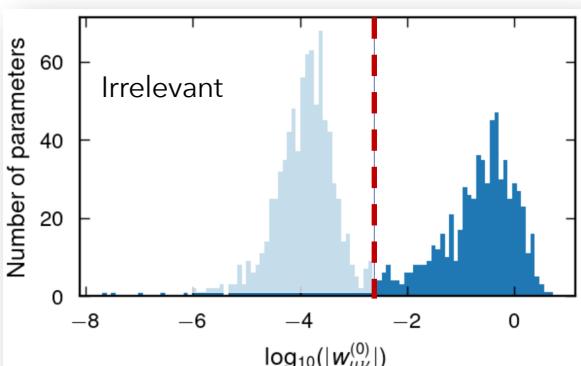
$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= \lambda_{\text{energy}} \cdot \text{RMSE}(\text{energy}) \\ &+ \lambda_{\text{forces}} \cdot \text{RMSE}(\text{forces}) \\ &+ \lambda_{\text{virial}} \cdot \text{RMSE}(\text{virial}) \\ &+ \lambda_1 \cdot \|\mathbf{w}\|_1 + \lambda_2 \cdot \|\mathbf{w}\|_2 \end{aligned}$$

Regularization  
 → stability  
 → data efficiency

$$\begin{aligned} E &= \sum_i E_i \\ E_i &= \sum_{\mu}^{\text{neu}} w_{\mu}^{(1)} \tanh \left( \sum_{\nu}^{\text{des}} w_{\mu\nu}^{(0)} q_{\nu}^{(i)} - b_{\mu}^{(0)} \right) - b^{(1)} \end{aligned}$$

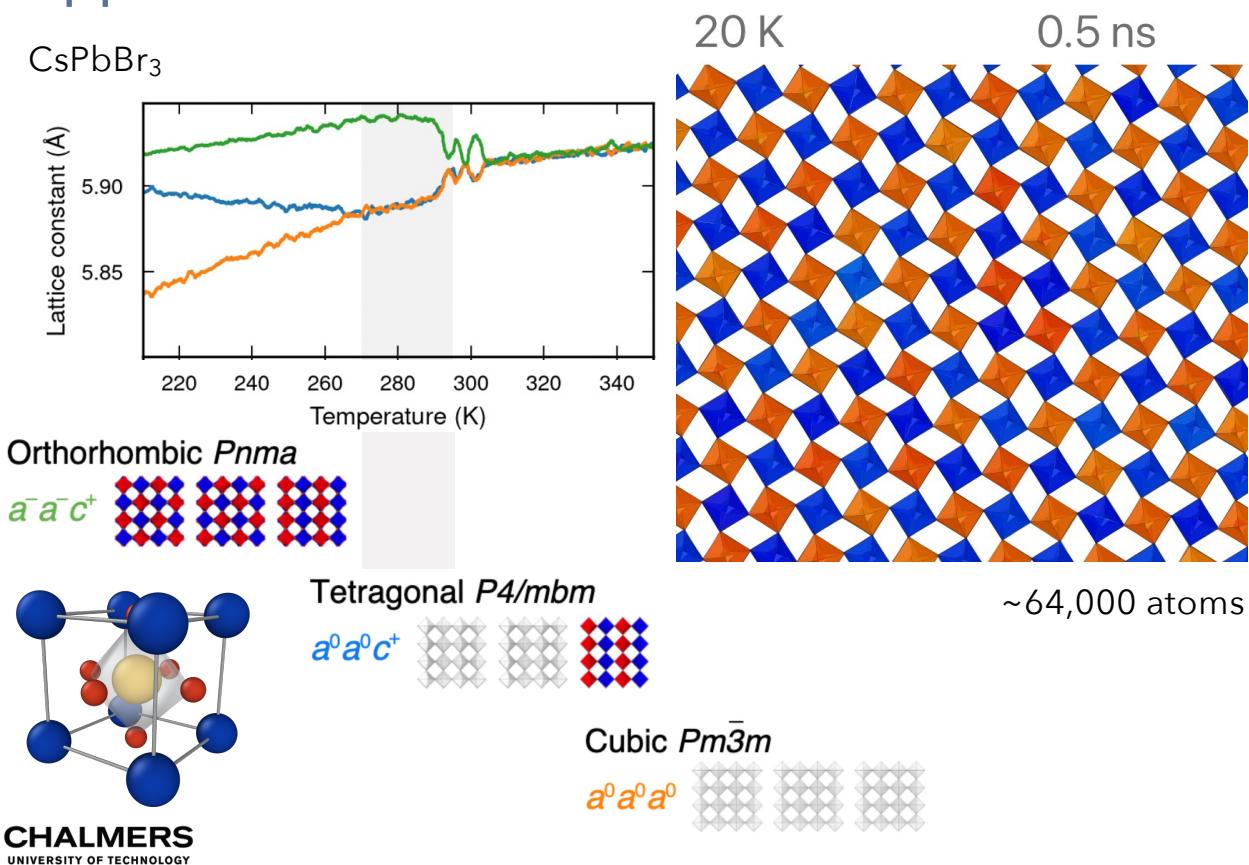
Optimization

separable natural evolution strategy



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

# Applications: Atomic scale models



9

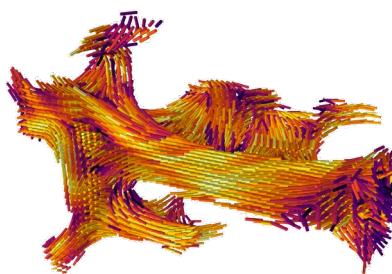
## Practical challenges

(Very) large **data** sets  
Incomplete data sets

(Very) large **parameter** spaces

Errors / **noise**

Potentially functional forms  
that are numerically problematic ("sloppy models")



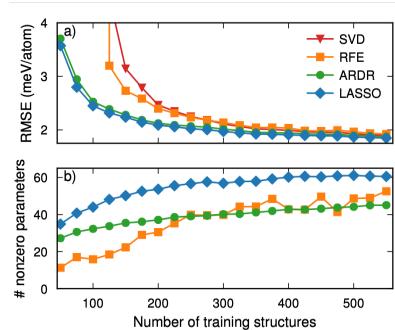
### Consequences

- Direct MCMC often impractical
- No general closed solutions for non-linear problems
- Model design and problem formulation important

## Follow-up questions (1)

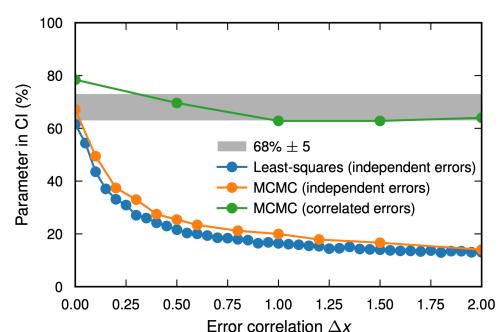
- How do we find out which parameters matter?  
→ sparse models via regularization
- How do we do this efficiently while accounting meaningfully for prior information?  
→ loss function

$$\mathcal{L} = \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} + \lambda_1 \|\boldsymbol{\theta}\|_1 + \lambda_2 \|\boldsymbol{\theta}\|_2^2 + \lambda_3 \|\nabla^2 \boldsymbol{\theta}\|_2^2 + \dots$$



- How we handle error correlation?  
→ "Sigma"

$$\mathcal{L} = \boldsymbol{\varepsilon}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\varepsilon} + \dots$$

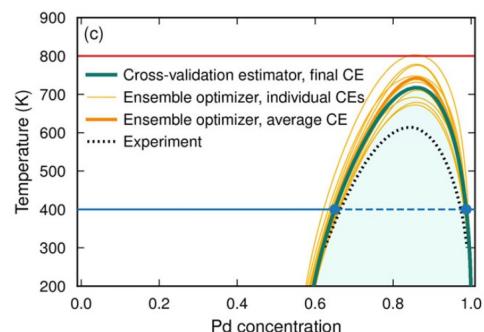


CHALMERS  
UNIVERSITY OF TECHNOLOGY

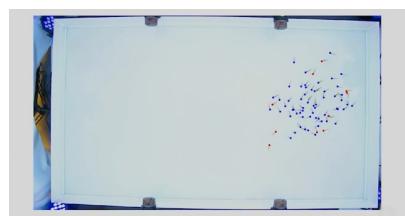
11

## Follow-up questions (2)

- How can we quantify model uncertainty?  
→ Sensitivity analysis and error propagation  
(e.g., bagging, committee/ensemble models, Mahalanobis norm),
- How do we acquire data?  
→ active learning



- How do find solutions in complex and often high-dimensional parameter spaces?  
→ global optimization



CHALMERS  
UNIVERSITY OF TECHNOLOGY

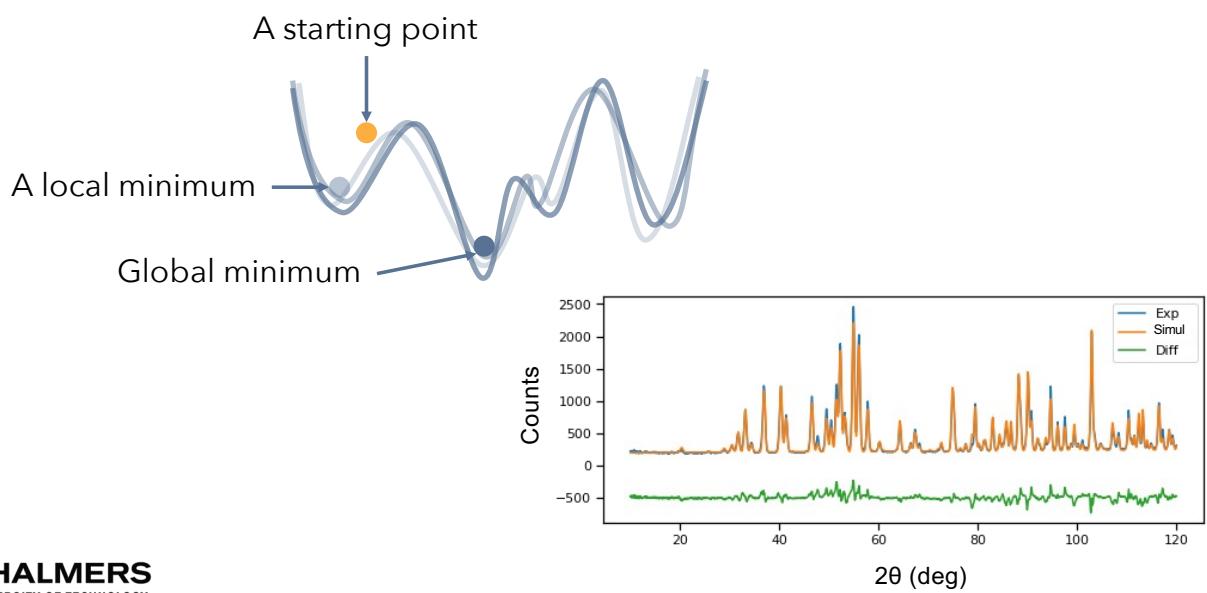
12

# Reconsider: What are the benefits of a Bayesian perspective?

13

## Problem statement

- |                                    |                                |
|------------------------------------|--------------------------------|
| 1. What do we optimize?            | → linear vs non-linear models  |
| 2. How do we optimize?             | → local vs global optimization |
| 3. What about errors?              | → error correlation            |
| 4. How do we convey prior insight? | → priors and covariances       |



# A motivating example

```
import numpy as np
from scipy.optimize import curve_fit
```

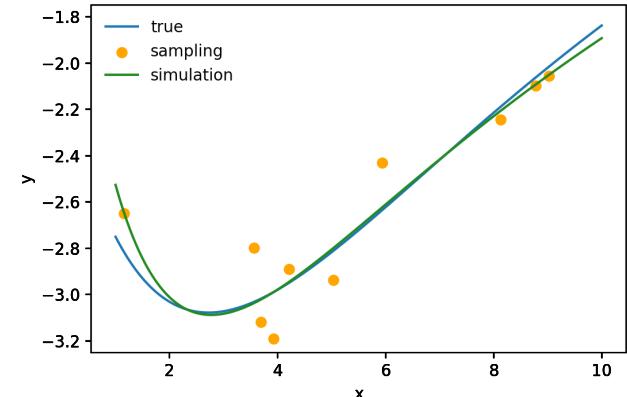
Truth

4 parameters

```
def my_func(x, a1, a2, t1, t2):
    return a1 * np.exp(-x / t1) \
        + a2 * np.exp(-(x - 0.1) / t2)

params_true = (3.0, -5.0, 2.0, 10.0)
```

$$y = \textcolor{red}{a}_1 \exp\left(-\frac{x}{t_1}\right) + \textcolor{red}{a}_2 \exp\left(-\frac{x - 0.1}{t_2}\right)$$



Sampling

```
xs_sampled = 9 * np.random.random_sample(size=10) + 1
ys_sampled = my_func(xs_sampled, *params_true)
ys_sampled += np.random.normal(scale=0.1, size=len(xs_sampled))
```

10 data points

add noise

Simulation/inference/training

```
params_simulated, params_cov = curve_fit(my_func, xs_sampled, ys_sampled)
```

Least-squares

**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

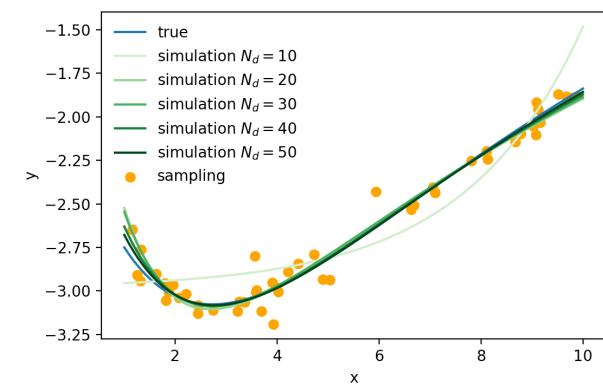
15

# A motivating example

- Adding more data does not necessarily imply convergence
- How can uncertainty be assessed?

Covariance

```
params_simulated, params_cov = curve_fit(
```

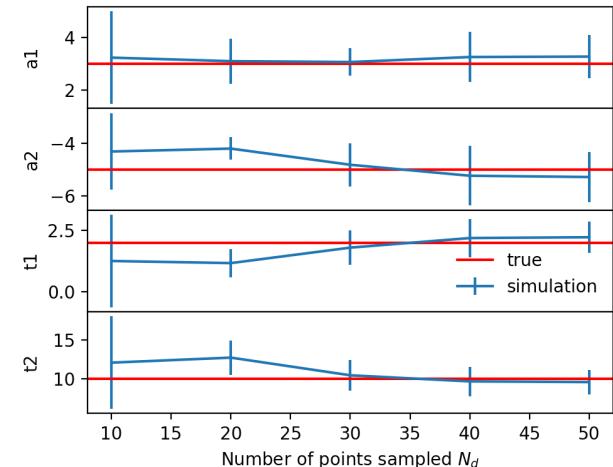


Parameter estimates with uncertainties

```
params_cov
```

```
array([[ 0.059, -0.074,  0.052, -0.204],
       [-0.074,  0.176, -0.159,  0.461],
       [ 0.052, -0.159,  0.152, -0.408],
       [-0.204,  0.461, -0.408,  1.222]])
```

```
params_errors = np.sqrt(np.diag(params_cov))
```

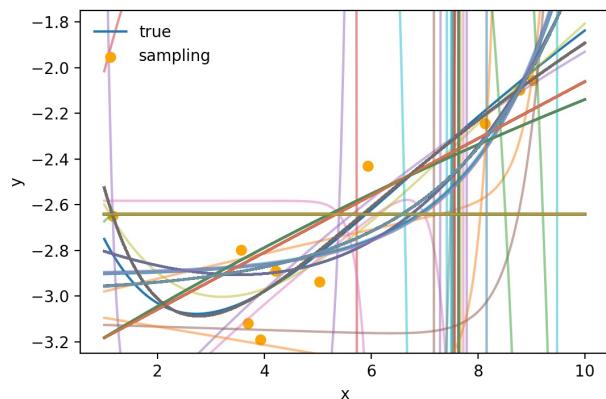


**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

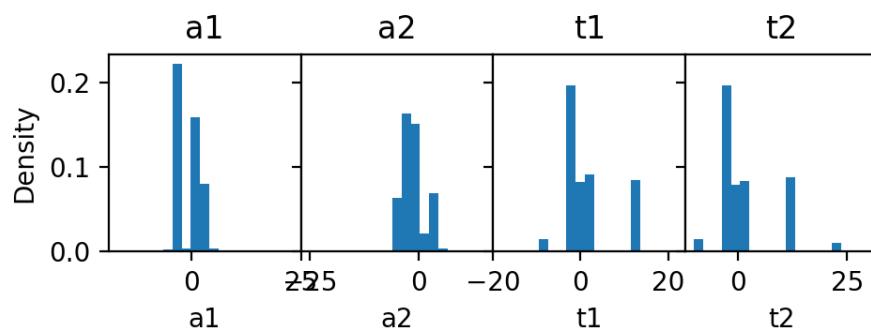
16

# A motivating example

Results can be very sensitive to the initial conditions



Distribution of "optimal" parameters



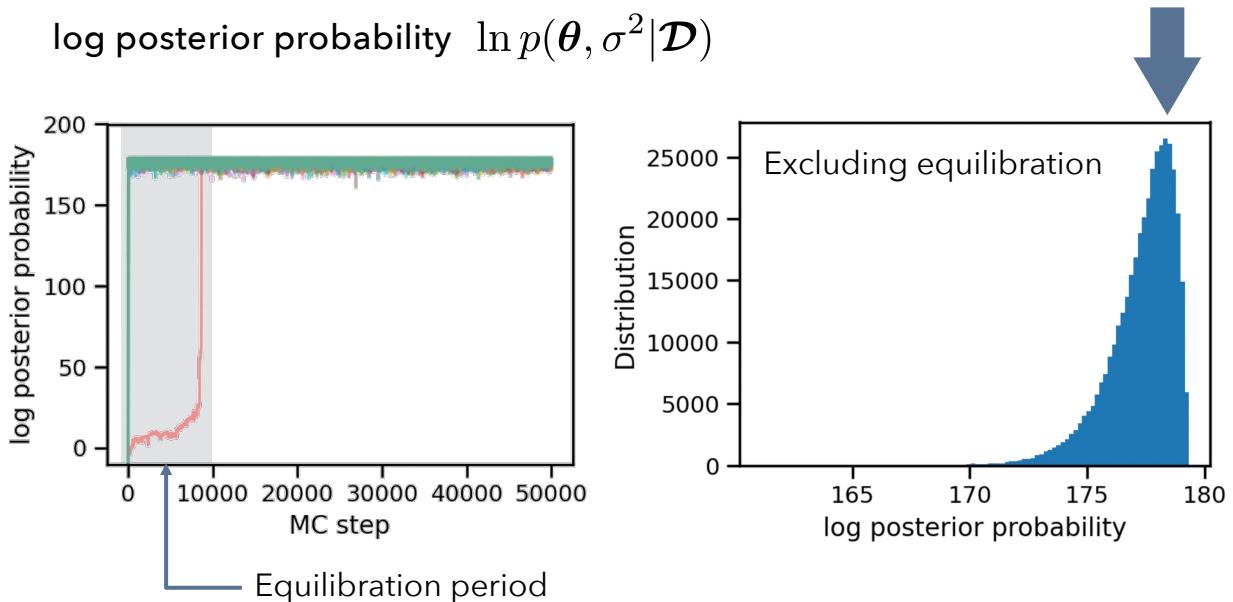
CHALMERS  
UNIVERSITY OF TECHNOLOGY

17

## Monte Carlo sampling

- Sampling of posterior via Markov chain-Monte Carlo simulations

Not a Gaussian!



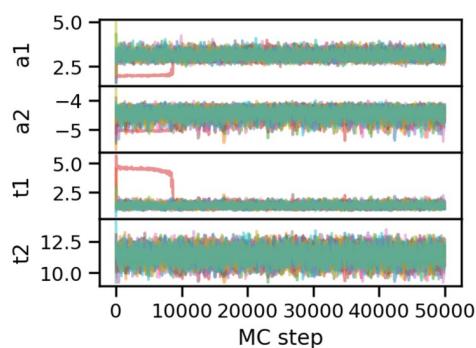
CHALMERS  
UNIVERSITY OF TECHNOLOGY

10

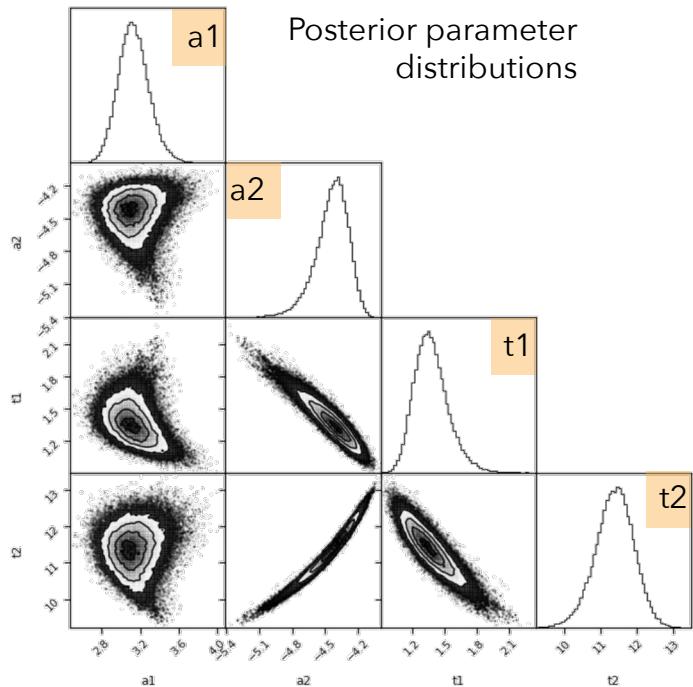
# Posterior parameter estimates



MCMC trajectories



Corner plot



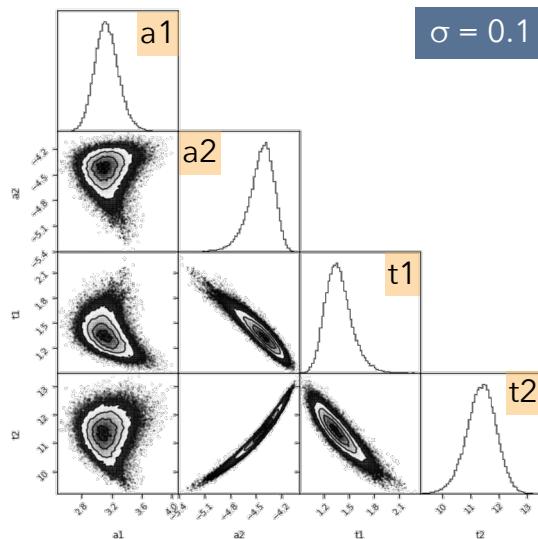
Posterior parameter distributions

How do error estimates  
in standard least-squares  
(curve\_fit) work?

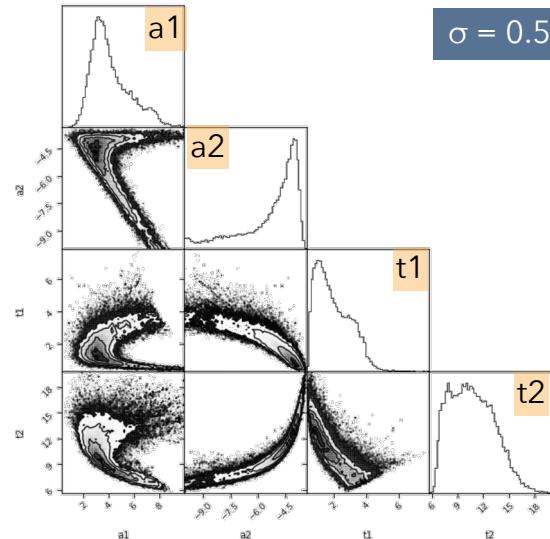
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

19

# Posterior parameter estimates



$\sigma = 0.1$



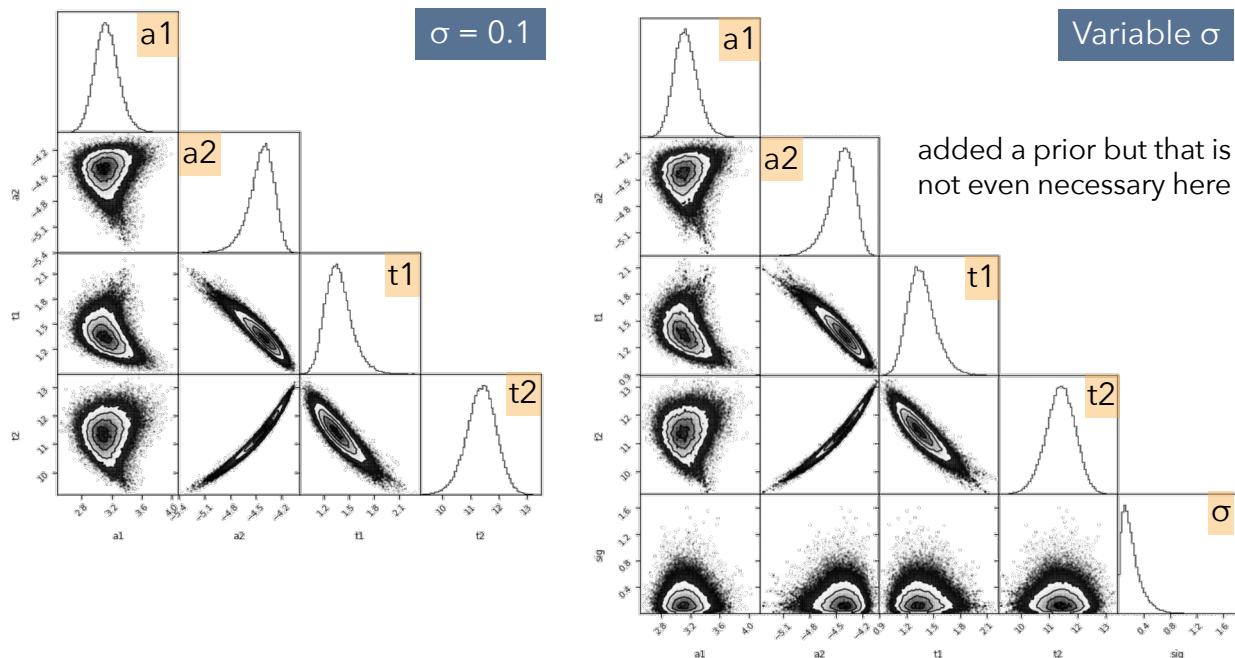
$\sigma = 0.5$

What is the “correct”  $\sigma$ ? → Sample it!

**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

20

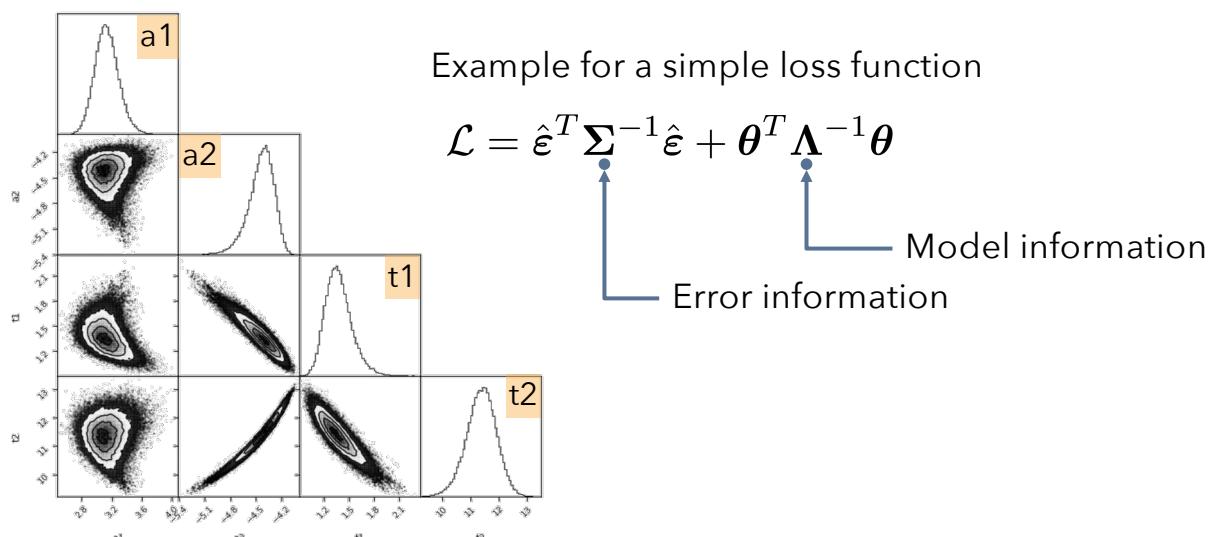
# Posterior parameter estimates



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

21

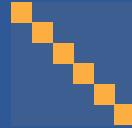
## What do we gain from a Bayesian perspective?



- Deeper understanding of correlations
- A probabilistic view that we can apply for interpretation

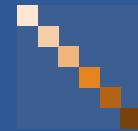
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

22



$$\Lambda = \lambda^2 \mathbf{1}$$

## Ridge regression and beyond: Parameter specific priors



$$\Lambda = \text{diag}$$

23

## Ridge and Bayesian ridge regression

Ridge regression  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | 0, \tau^2)$

Bayesian ridge regression  $p(\mathcal{D}; M(\boldsymbol{\theta}), \alpha) = \mathcal{N}(\mathcal{D} | M(\boldsymbol{\theta}), \alpha)$   
 $p(\boldsymbol{\theta} | \lambda) = \mathcal{N}(\boldsymbol{\theta} | 0, \lambda^{-1} \mathbf{I})$

$\alpha$  and  $\lambda$  from Gamma distributions  
estimated jointly during fitting by maximizing  
the log marginal likelihood (LML)

Recall

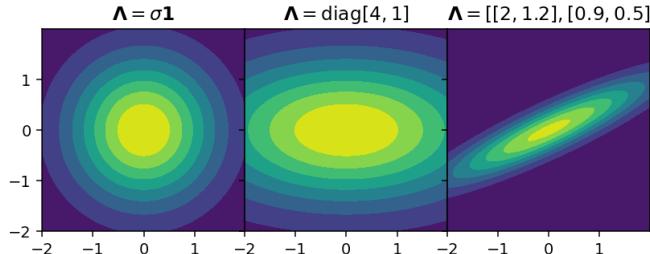
$$\mathcal{L} = \hat{\boldsymbol{\varepsilon}}^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\varepsilon}} + \boldsymbol{\theta}^T \boldsymbol{\Lambda}^{-1} \boldsymbol{\theta}$$

↑      ↑  
  Error information    Model information

# Automatic relevance detection (ARD) regression

$$p(\boldsymbol{\theta}|\Lambda) = \mathcal{N}(\boldsymbol{\theta}|0, \Lambda^{-1})$$
$$\text{diag}\Lambda = [\lambda_1, \lambda_2 \dots \lambda_{N_p}]$$

Now: Gamma distribution  
for each parameter



Followed by a pruning step  
where parameters for which the  
statistical significance is below a  
threshold are set to zero  
→ Feature selection (next)

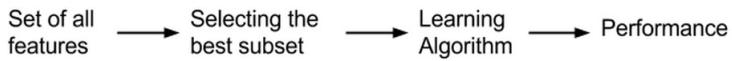
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

25

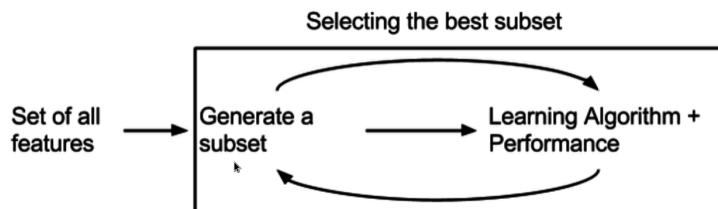
**Feature selection:  
Making models sparse**

# (Automated) feature selection

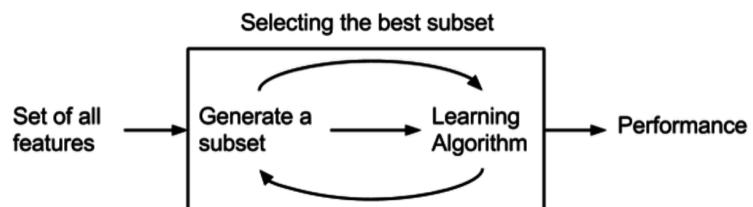
Filter-type



Embedded-type



Wrapper-type (RFE, genetic algorithms, ant colony optimization, ...)



# (Automated) feature selection

## "Embedded"-type

Algorithms that "take advantage of their own variable selection"

Automatic relevance detection (ARD) regression

$$p(\boldsymbol{\theta}|\boldsymbol{\Lambda}) = \mathcal{N}(\boldsymbol{\theta}|0, \boldsymbol{\Lambda}^{-1})$$
$$\text{diag}\boldsymbol{\Lambda} = [\lambda_1, \lambda_2 \dots \lambda_{N_p}]$$

Followed by pruning step

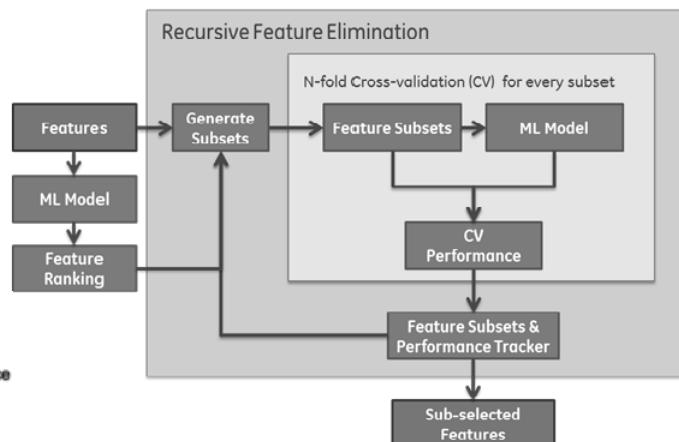
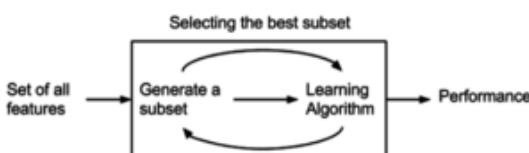
$$\hat{\theta}_i = \begin{cases} \theta_i & \text{if } \lambda_i < \lambda_{\text{threshold}} \\ 0 & \text{else} \end{cases}$$

Principle can be similarly applied to LASSO

# (Automated) feature selection

## Wrapper-type

Algorithms that wrap around another regression technique  
 → Example: Recursive feature elimination (RFE)



```

from sklearn.feature_selection import RFE, RFECV
from sklearn.linear_model import LinearRegression

X, E = sc.get_fit_data()
est = LinearRegression(...)
rfe = RFE(estimator=est, n_features_to_select=...)
rfe.fit(X, E)
    
```

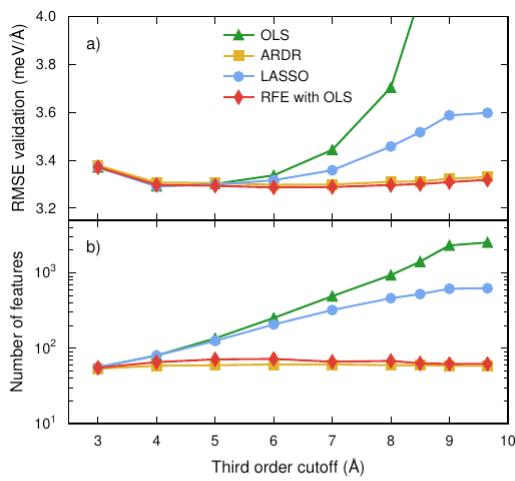
hyperparameter



Figures from wiki commons and doi:10.1109/EMBC.2016.7591621

29

## Ex: Thermal conductivity from force constants



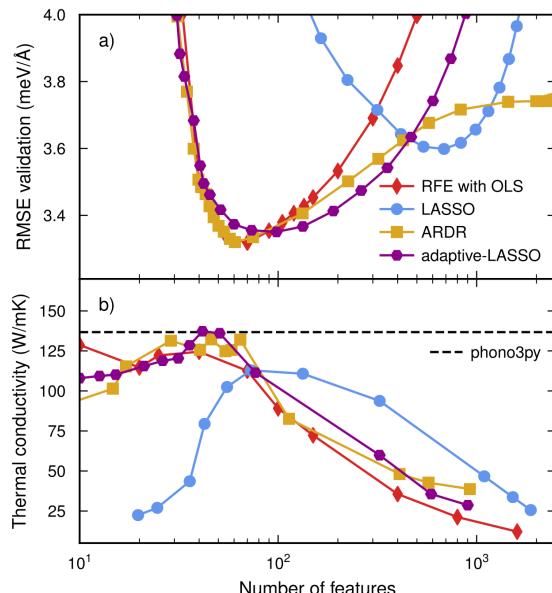
Nick Drachman  
 @ndrach1

So these "higher order terms", are they in the room with us right now?



Sparsity often implies  
 "more physical" models

Often improves stability



# Dealing with outliers: Robust regression

31

## Dealing with non-normal iid errors

**OLS, ridge, LASSO etc:**  
errors given by normal distribution (and iid)

$$p(\mathcal{D}; M(\boldsymbol{\theta}), \alpha) = \mathcal{N}(\mathcal{D}|M(\boldsymbol{\theta}), \alpha)$$

What happens when you have outliers or “heavy tails”?  
→ notebooks 1 and 2

**Robust regression:**  
assume heavy-tailed distribution

$$p(\mathcal{D}; M(\boldsymbol{\theta}), b) = \exp\left(-\frac{1}{b}|\mathcal{D} - M(\boldsymbol{\theta})|\right) \quad \text{e.g., Laplace}$$

# The Huber regressor

**Robust regression:**

assume heavy-tailed distribution

$$p(\mathcal{D}; M(\boldsymbol{\theta}), b) = \exp\left(-\frac{1}{b}|\mathcal{D} - M(\boldsymbol{\theta})|\right) \quad \text{e.g., Laplace}$$

Problem: Non-differentiable loss function

$$\mathcal{L} = \|\mathcal{D} - M(\boldsymbol{\theta})\|_1 + \dots$$

Solution: linear programming or **Huber regressor**

$$\mathcal{L} = H_\epsilon(\mathcal{D} - M(\boldsymbol{\theta})) + \dots$$

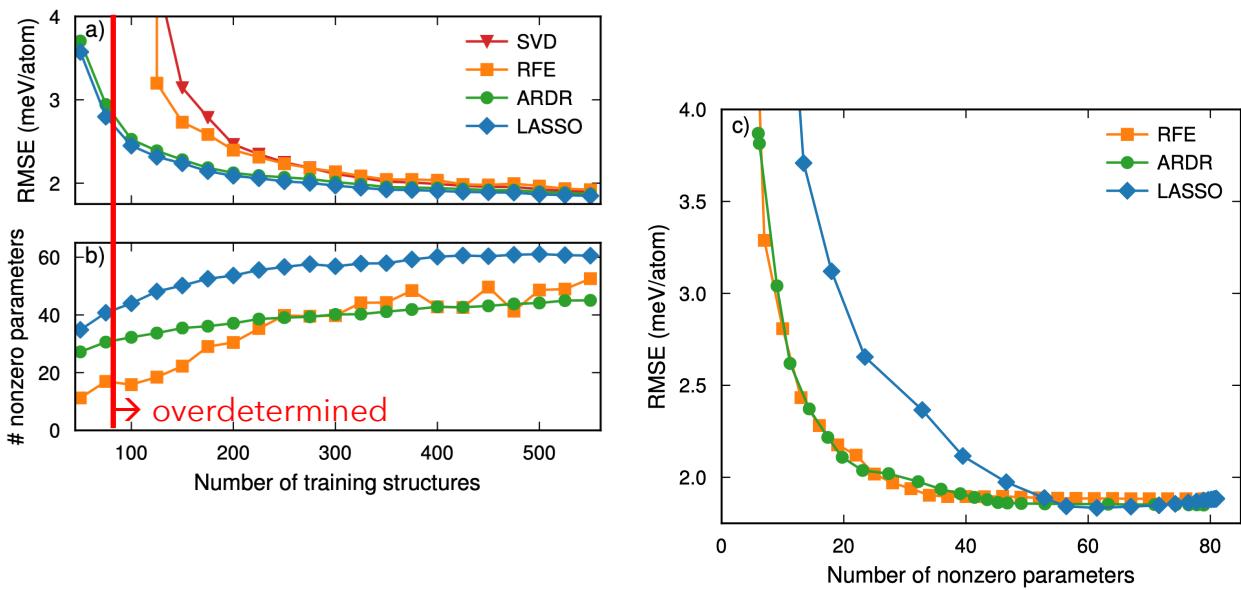
$$H_\epsilon(z) = \begin{cases} z^2 & \text{for } |z| \leq \epsilon \\ 2\epsilon|z| - \epsilon^2 & \text{otherwise} \end{cases}$$



33

Method performance comparison

# Example: Model for Ag-Pd alloy



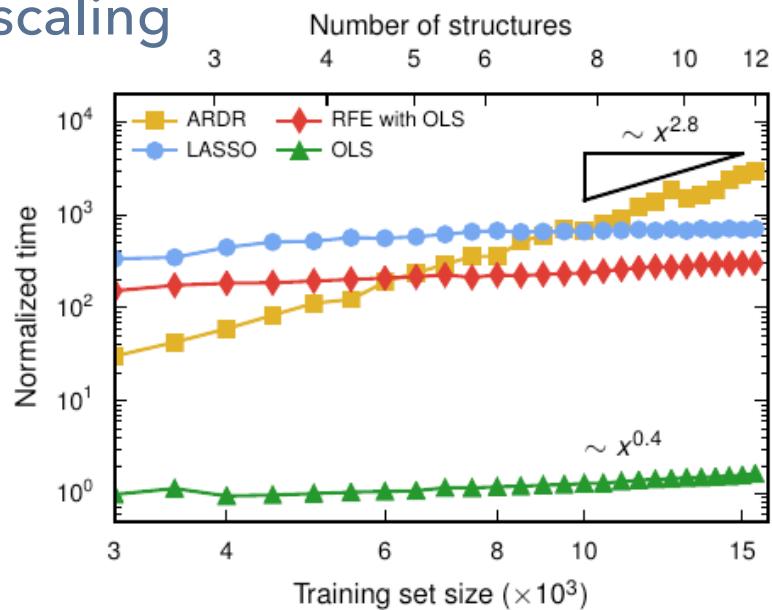
Cluster expansion  
→ Wednesday lecture

ARDR: fast convergence and sparse solution  
LASSO: more false-positives  
RFE: requires more structures

**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

35

## Computational scaling



Cost becomes significant for large problems  
ARDR (Bayesian) becomes too costly due to scaling  
RFE or OLS with manual selection often the most performant

Force constant expansion  
→ lecture notes

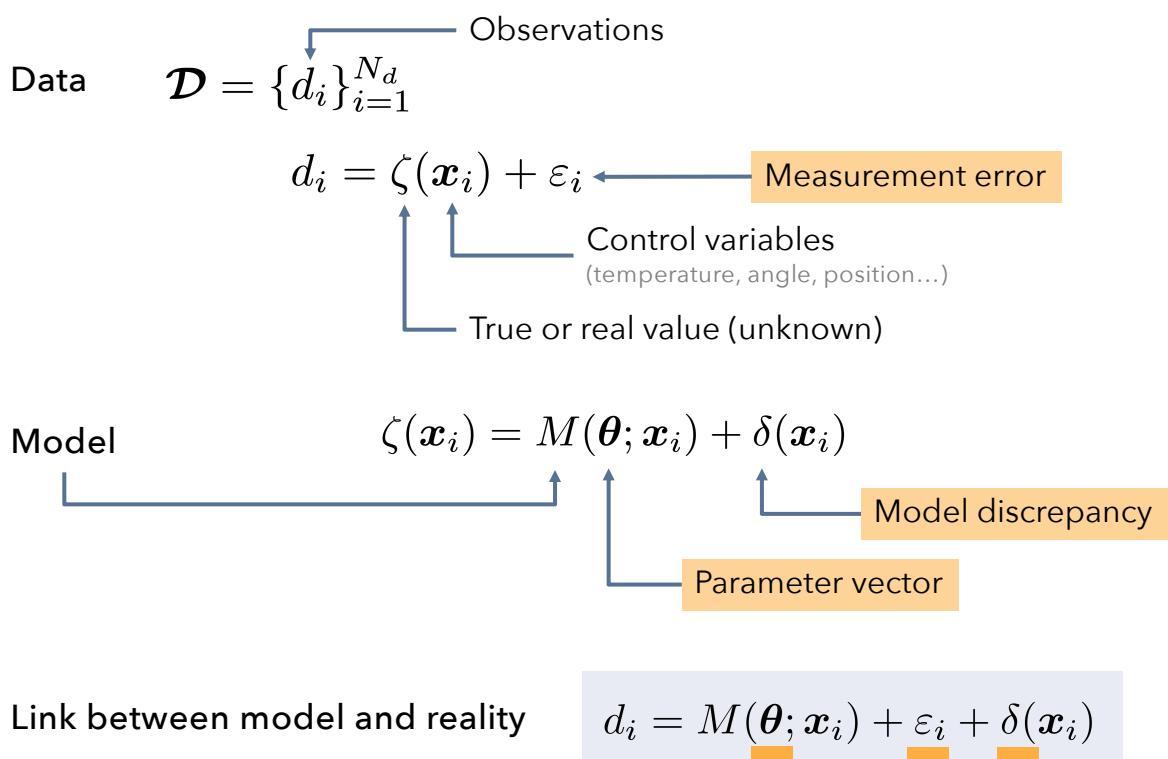
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

Fransson *et al.*, npj Comp. Mater. **6**, 135 (2020)

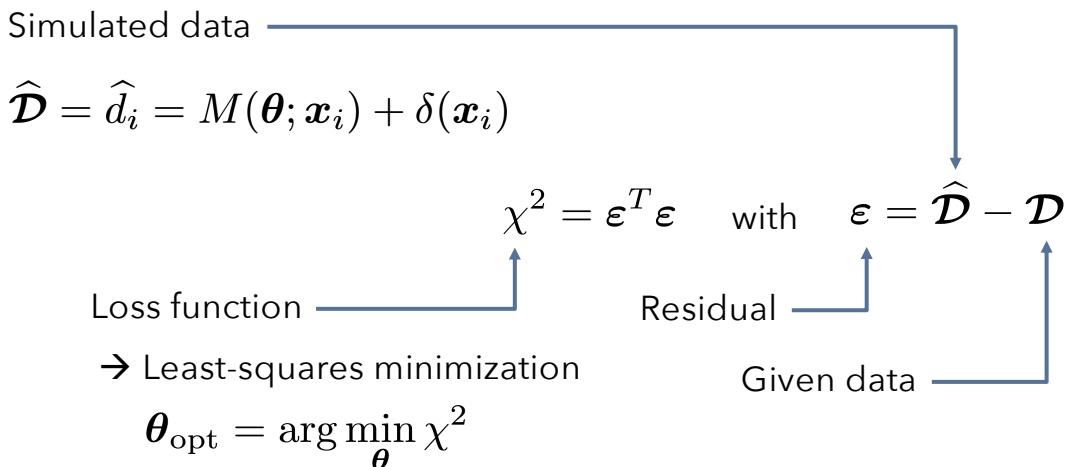
## Error correlation

37

## Simulating reality



# What did we do before?



1. We did not *consciously* model the error
2. We did not use any prior knowledge regarding the error

## Modeling more complex error patterns

Generalize the variance = go beyond i.i.d. errors

$$\frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{\sigma^2} \rightarrow \boldsymbol{\varepsilon}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\varepsilon}$$

Covariance matrix

$$p(\mathcal{D}|\boldsymbol{\theta}, \sigma^2) = \left( \frac{1}{2\pi\sigma^2} \right)^{N_d/2} \exp \left\{ -\frac{1}{2} \boldsymbol{\varepsilon}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\varepsilon} \right\}$$



$\boldsymbol{\Sigma} = \sigma^2 \mathbf{1}$   
i.i.d.



$\boldsymbol{\Sigma} = \text{diag}$   
e.g., error varies with  $2\theta$  or signal intensity

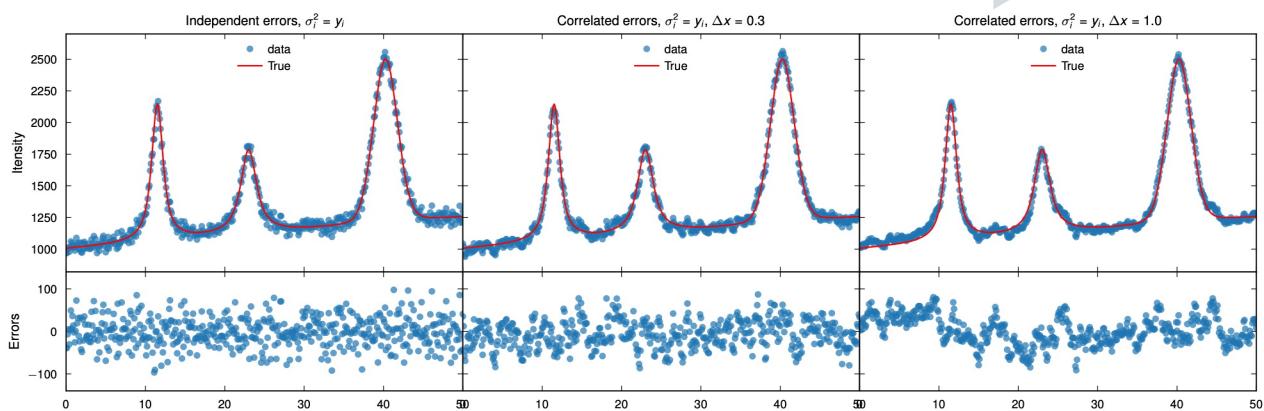


$\boldsymbol{\Sigma} = \text{full}$   
e.g., correlated errors

# Modeling more complex error patterns

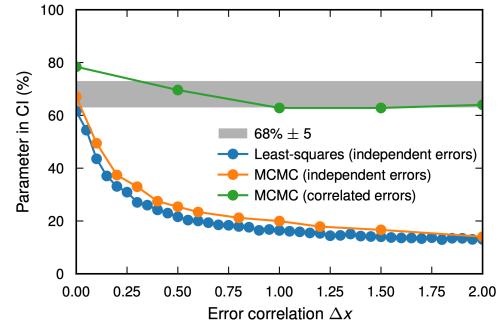
A demonstration

Increasing correlation



$$\Sigma_{ij} = \sigma^2 \exp(-|x_i - x_j|/\Delta x)$$

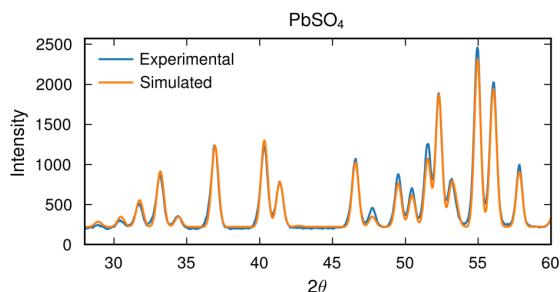
If error correlations are untreated  
the parameter uncertainty is  
dramatically underestimated



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

41

## Example: Refining diffraction data



$y_{\text{exp}}$  Measured intensity

$y_{\text{sim}}$  Simulated intensity

$\mathbf{x}$  Instrumental parameters,  
structural parameters etc.

### Minimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \chi^2(\mathbf{x})$$

$$\chi^2(\mathbf{x}) = \sum_i \frac{1}{\sigma^2} \frac{(y_{i,\text{sim}} - y_{i,\text{exp}})^2}{y_{i,\text{exp}}^2}$$

Likelihood from assumption  
of identical normal errors

$$y_{i,\text{exp}} = y_{i,\text{sim}} + \varepsilon_i$$

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2 y_{i,\text{exp}}^2)$$

**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

42

# MCMC vs Least-squares

$$p(\mathbf{x}) = \exp(-\chi^2(\mathbf{x}))$$

$$\chi^2(\mathbf{x}) = \sum_i \frac{1}{\sigma^2} \frac{(y_{i,\text{sim}} - y_{i,\text{exp}})^2}{y_{i,\text{exp}}^2}$$

## least squares (lmfit)

$$\begin{aligned} a &= 8.47701 \pm 0.00027 \text{ \AA} \\ b &= 5.39576 \pm 0.00019 \text{ \AA} \\ c &= 6.95584 \pm 0.00032 \text{ \AA} \end{aligned}$$

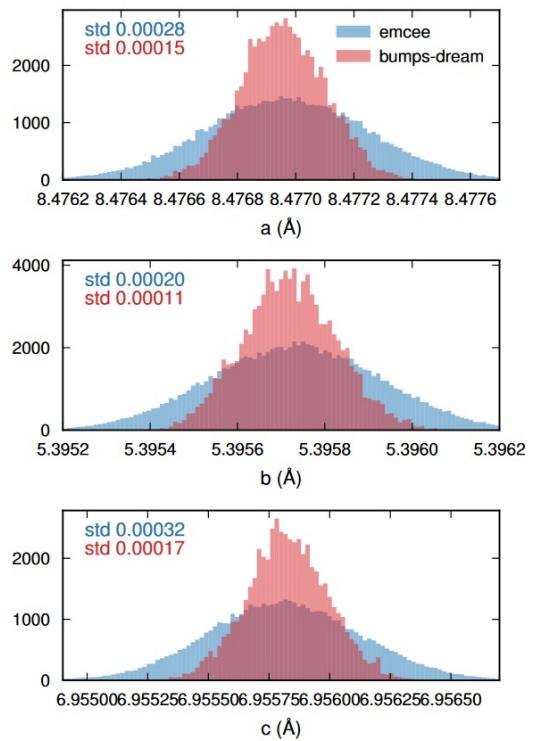
## MCMC (emcee) $\sigma \approx 2.0$

$$\begin{aligned} a &= 8.47697 \pm 0.00028 \text{ \AA} \\ b &= 5.39564 \pm 0.00019 \text{ \AA} \\ c &= 6.95585 \pm 0.00031 \text{ \AA} \end{aligned}$$

## MCMC (bumps, dream; no correlation)

$$\begin{aligned} a &= 8.47699 \pm 0.00015 \text{ \AA} \\ b &= 5.39571 \pm 0.00011 \text{ \AA} \\ c &= 6.95584 \pm 0.00017 \text{ \AA} \end{aligned}$$

Smaller  $\sigma \rightarrow$  Smaller errors  
More data  $\rightarrow$  Smaller errors

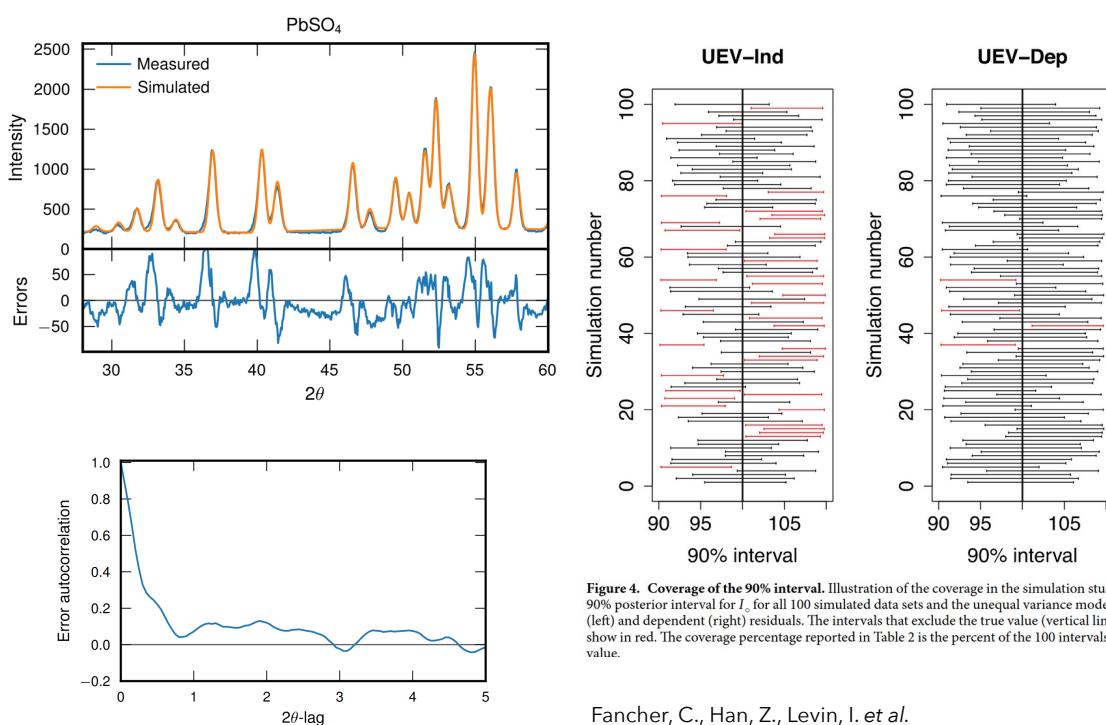


**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

43

# Error correlations

$\varepsilon_i$  correlates with  $\varepsilon_j$



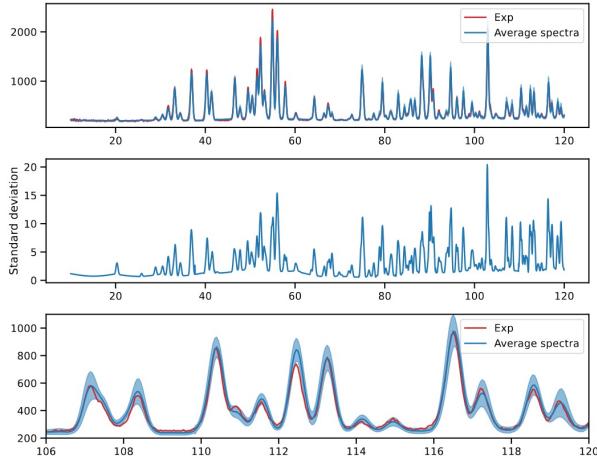
**Figure 4.** Coverage of the 90% interval. Illustration of the coverage in the simulation study. The plot shows the 90% posterior interval for  $I_c$  for all 100 simulated data sets and the unequal variance models with independent (left) and dependent (right) residuals. The intervals that exclude the true value (vertical line at  $I_c = 100$ ) are shown in red. The coverage percentage reported in Table 2 is the percent of the 100 intervals that include the true value.

Fancher, C., Han, Z., Levin, I. et al.  
Use of Bayesian Inference in Crystallographic Structure Refinement via Full Diffraction Profile Analysis  
*Sci Rep* **6**, 31625 (2016); <https://doi.org/10.1038/srep31625>

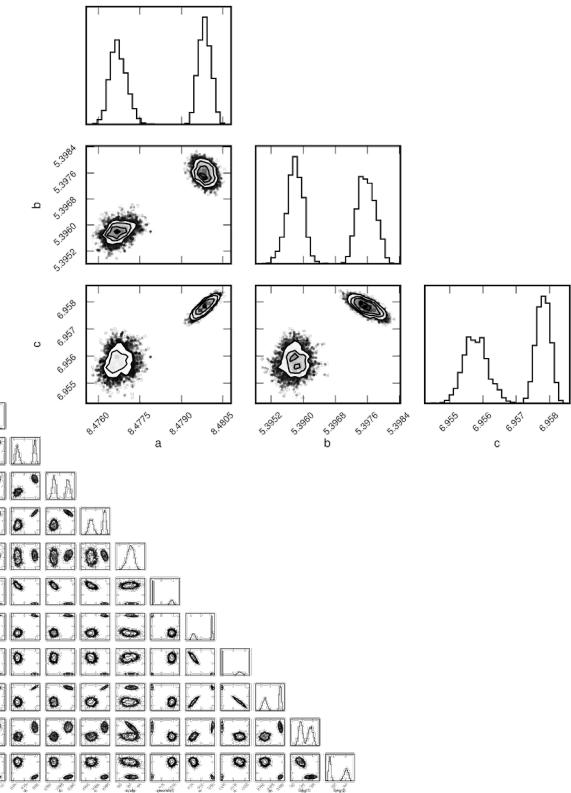
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

11

# Ensemble of spectra



Multiple minima for independent MCMC walkers



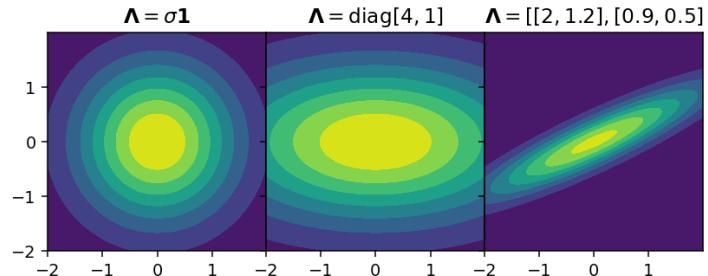
**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

45

## Take aways

$$\mathcal{L} = \hat{\epsilon}^T \Sigma^{-1} \hat{\epsilon} + \theta^T \Lambda^{-1} \theta$$

↑      ↑  
Model information    Error information



## Schematic of relation between methods

Method	Likelihood	Parameter prior	Hyperparameter prior
<b>OLS</b>			
Ridge			
Bayesian ridge			
<b>ARD</b>			
LASSO			
Robust regression			

## Error correlation

