



UNIVERSITY OF
CAMBRIDGE



Institute of
Computing for
Climate Science

Using Gaussian Processes to Uncover the Secrets of the Universe

Henry Moss @ MASCOTNUM 2024



UNIVERSITY OF
CAMBRIDGE



Institute of
Computing for
Climate Science

Using Gaussian Processes to Uncover the Secrets of the Universe

Stochastic Equation Discovery via Interpretable Additive Models

Henry Moss @ MASCOTNUM 2024

Why do
we want to
learn
symbolic
equations?

La presqu'île de Giens



Why do
we want to
learn
symbolic
equations?

La presqu'île de Giens



Morecambe



Why do
we want to
learn
symbolic
equations?

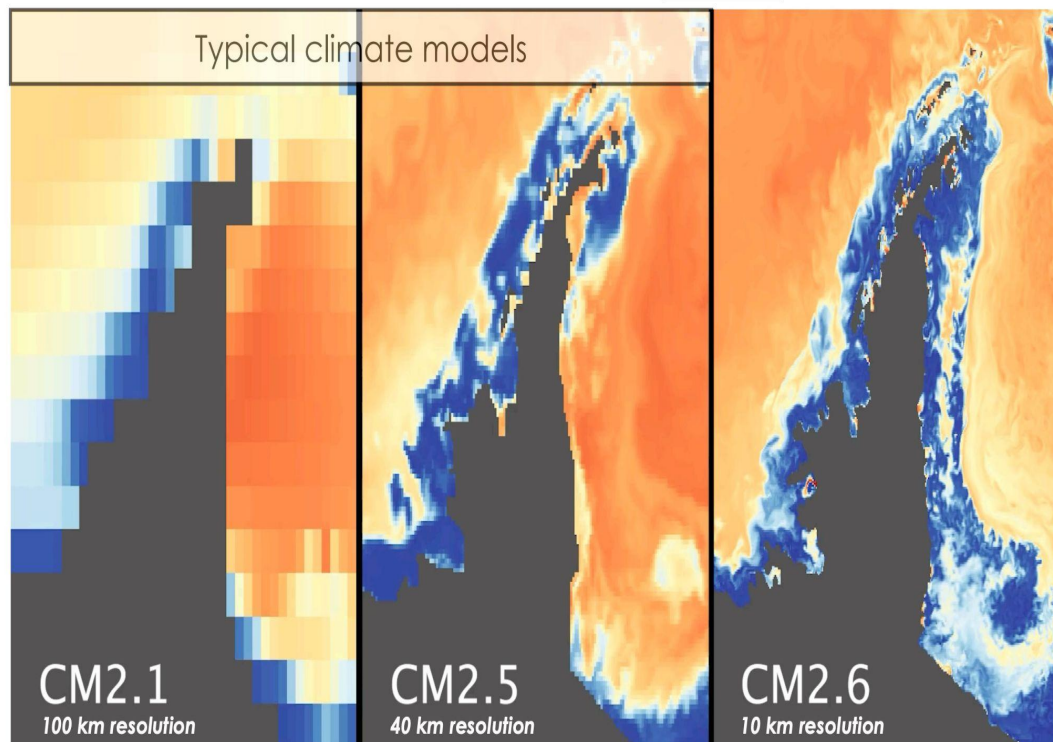


Why do
we want to
learn
symbolic
equations?



To learn parameterisations

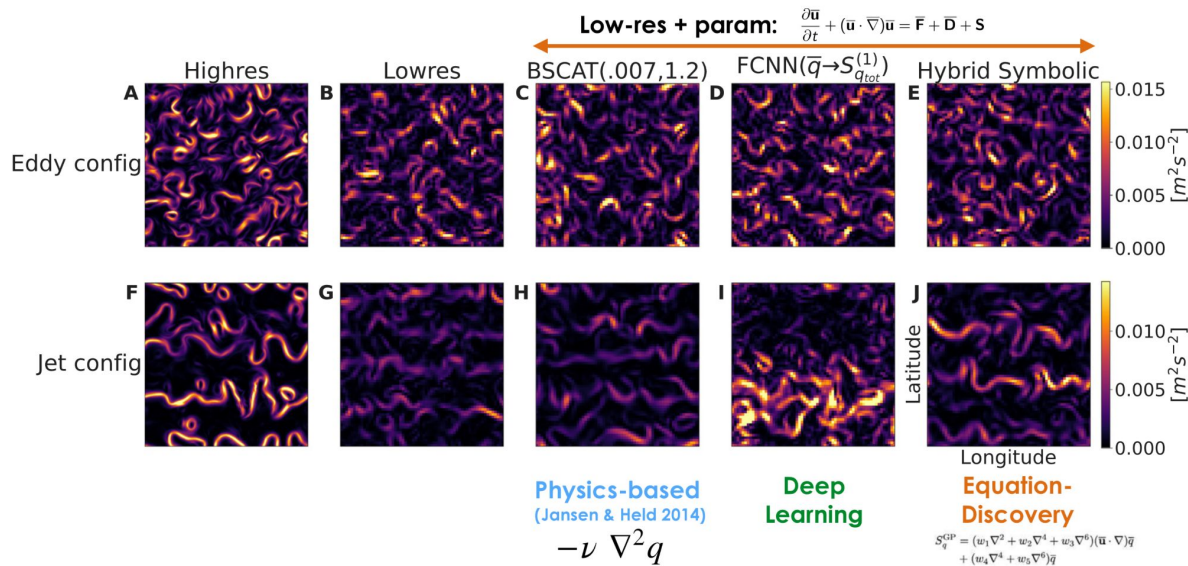
Why do we want to learn symbolic equations?



NOAA GFDL CM2 Suite; Animation from J. Busecke

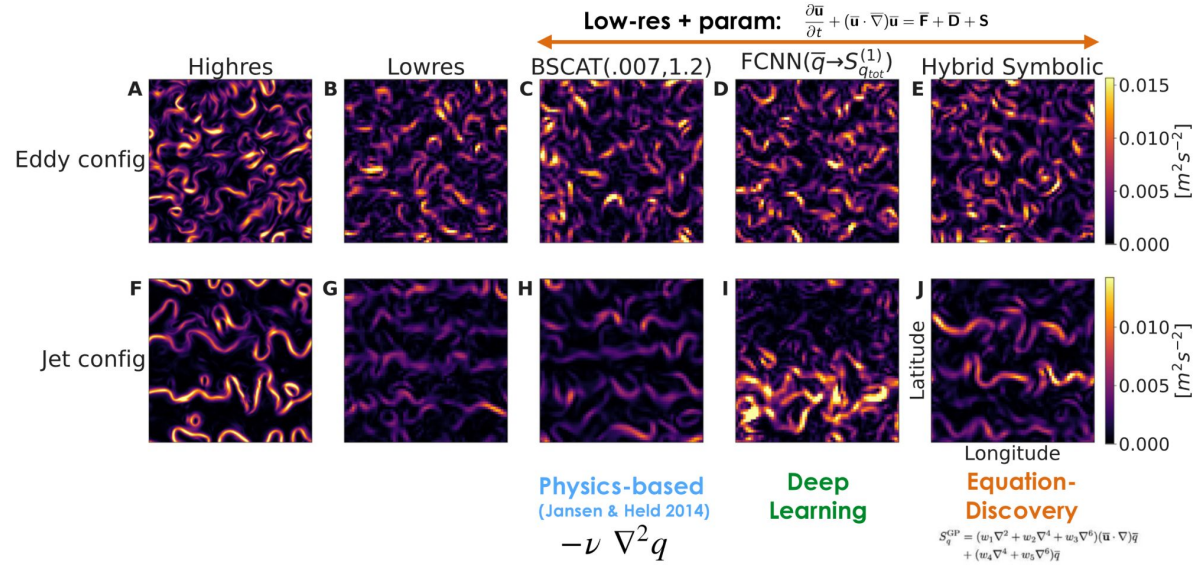
Why do we want to learn symbolic equations?

To learn parameterisations



To learn parameterisations

Why do we want to learn symbolic equations?



Zanna et al 2020

$$C_{\text{Sundqvist}} \stackrel{\text{def}}{=} 1 - \sqrt{\frac{\min\{\text{RH}, \text{RH}_{\text{sat}}\} - \text{RH}_{\text{sat}}}{\text{RH}_0 - \text{RH}_{\text{sat}}}} \quad C_{\text{Teixeira}} \stackrel{\text{def}}{=} \frac{Dq_c}{2q_s(1 - \widehat{\text{RH}})K} \left(-1 + \sqrt{1 + \frac{4q_s(1 - \widehat{\text{RH}})K}{Dq_c}} \right)$$

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = I_1(\text{RH}, T) + I_2(\partial_z \text{RH}) + I_3(q_c, q_i),$$

Grundner et al. 2023

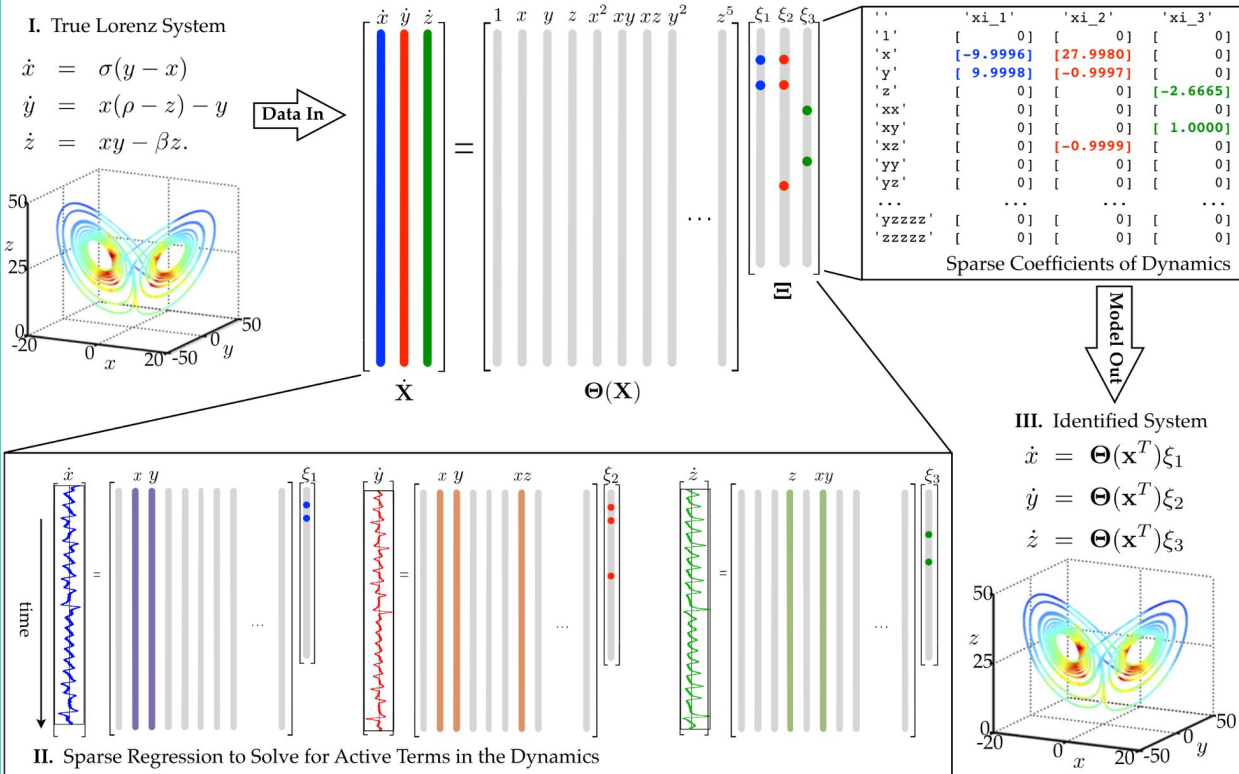
How can
we learn
symbolic
equations?

How can
we learn
symbolic
equations?

E.g. Sparse Identification of
Nonlinear Dynamics

How can we learn symbolic equations?

E.g. Sparse Identification of Nonlinear Dynamics

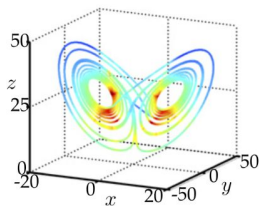


How can we learn symbolic equations?

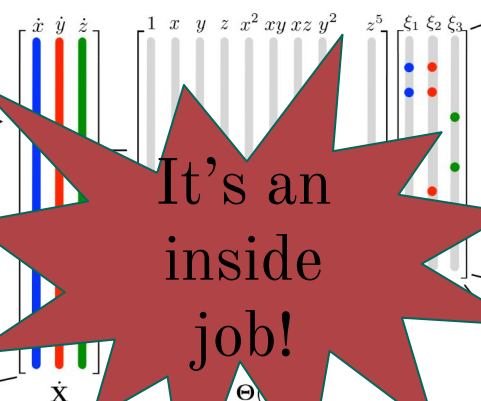
E.g. Sparse Identification of Nonlinear Dynamics

I. True Lorenz System

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= x(\rho - z) - y \\ \dot{z} &= xy - \beta z.\end{aligned}$$



Data In



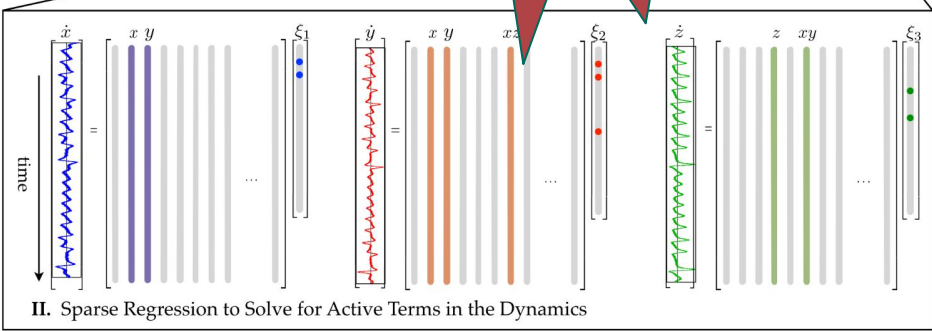
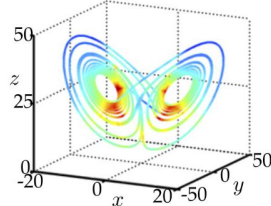
	'xi_1'	'xi_2'	'xi_3'
'1'	[0]	[0]	[0]
'x'	[-9.9996]	[27.9980]	[0]
'y'	[9.9998]	[-0.9997]	[0]
'z'	[0]	[0]	[-2.6665]
'xx'	[0]	[0]	[0]
'xy'	[0]	[0]	[1.0000]
'xz'	[0]	[-0.9999]	[0]
'yy'	[0]	[0]	[0]
'yz'	[0]	[0]	[0]
...
'yzzzz'	[0]	[0]	[0]
'zzzzz'	[0]	[0]	[0]

Sparse Coefficients of Dynamics

Model Out

III. Identified System

$$\begin{aligned}\dot{x} &= \Theta(\mathbf{x}^T)\xi_1 \\ \dot{y} &= \Theta(\mathbf{x}^T)\xi_2 \\ \dot{z} &= \Theta(\mathbf{x}^T)\xi_3\end{aligned}$$



II. Sparse Regression to Solve for Active Terms in the Dynamics


- Functional form in advance
- Correlated inputs
- Only uncertainty over params

What do we want?

- ML to **HELP** scientists discovery equations
- Learn **STOCHASTIC** equations


Lets use Gaussian processes!

—

A red, multi-pointed starburst graphic with a thin black outline, centered on the left side of the slide. It contains the text 'They can't handle large data volumes' in a black serif font.

They can't
handle large
data volumes


Let's use Gaussian processes!



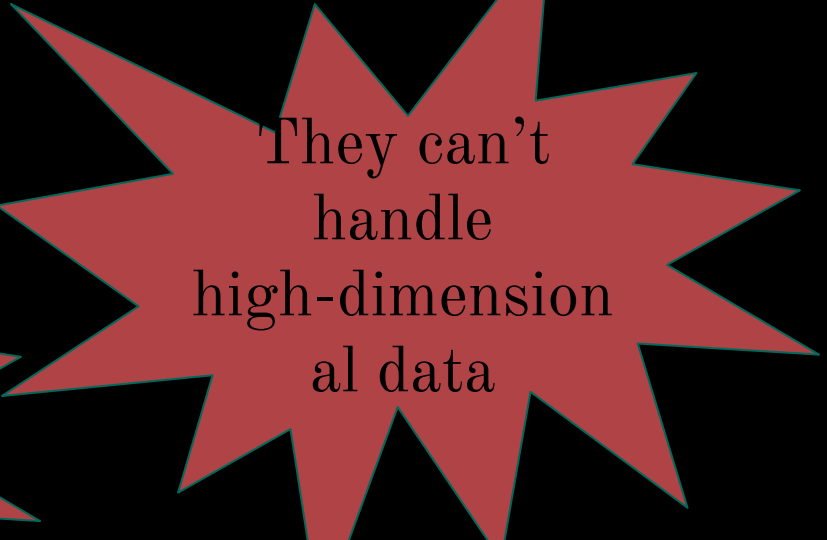
They can't
handle large
data volumes

Only for
Gaussian
data.....

Gaussian processes!



They can't
handle large
data volumes



They can't
handle
high-dimension
al data



Only for
Gaussian
data.....

Gaussian processes!

They can't
handle large
data volumes

They can't
handle
high-dimension
al data

Only for
Gaussian
data.....

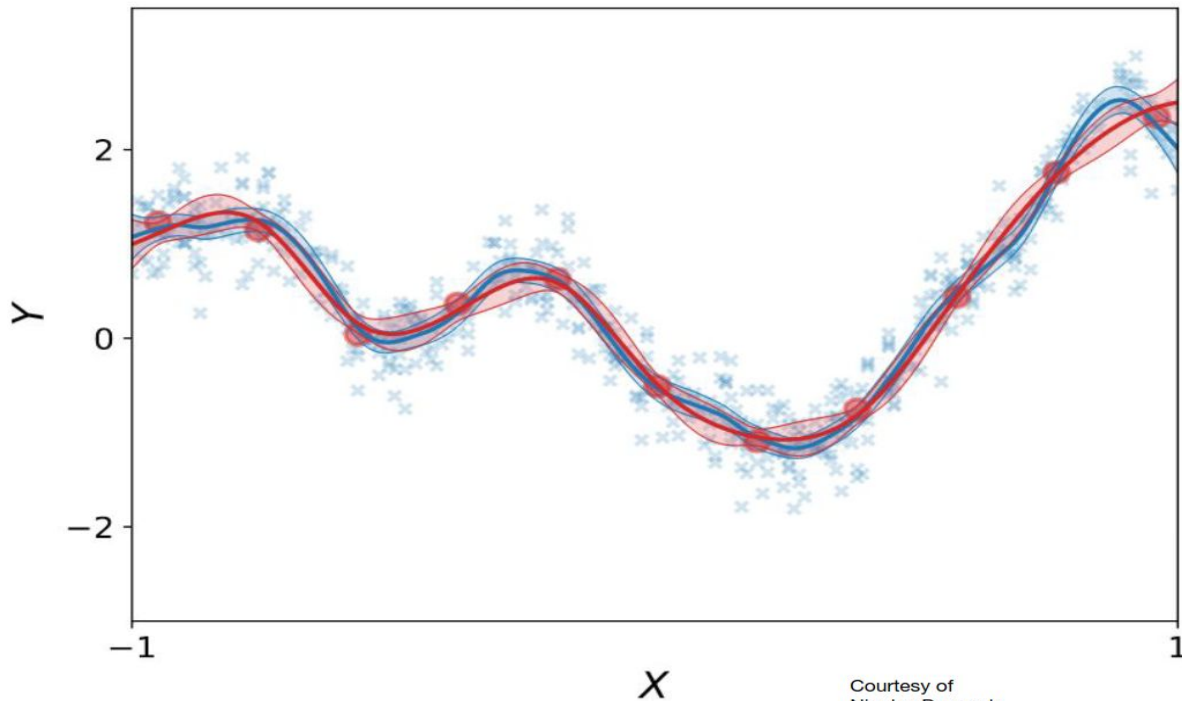
They are not
interpretable
(symbolic)

Linear Gaussian Processes!

GPs for big data?

- Use Sparse variational GP
- Replace with $M \ll N$

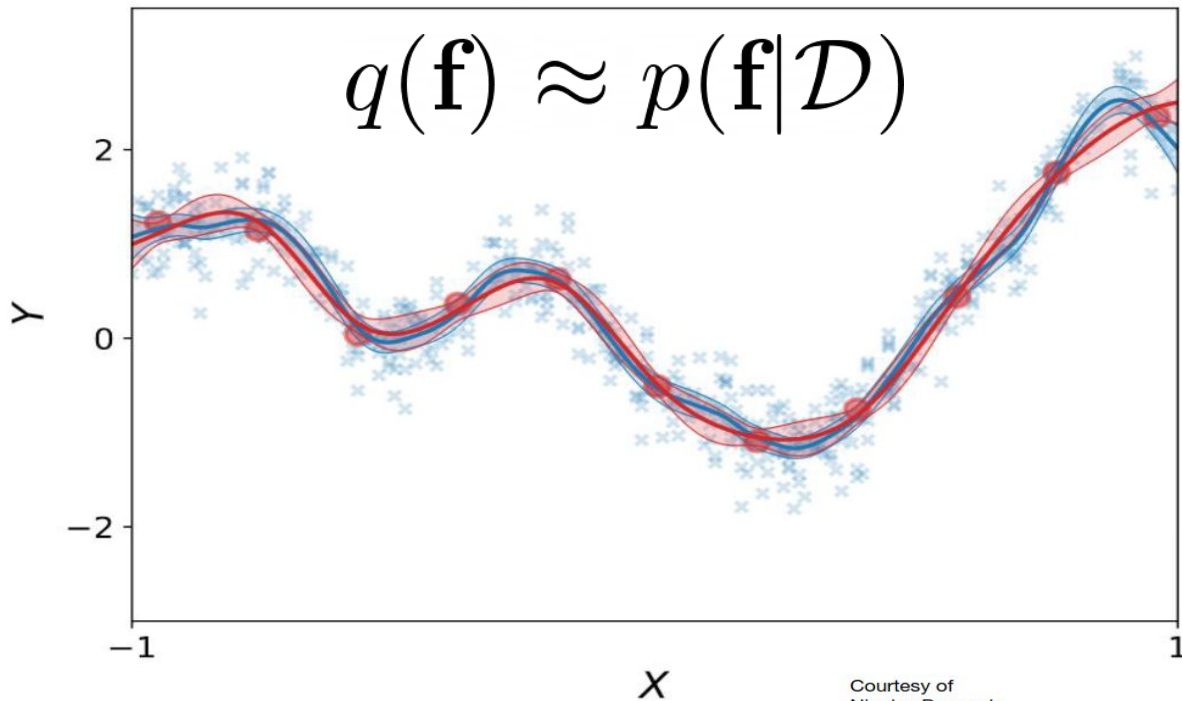
representative points



GPs for big data?

- Use Sparse variational GP
- Replace with $M \ll N$

representative points

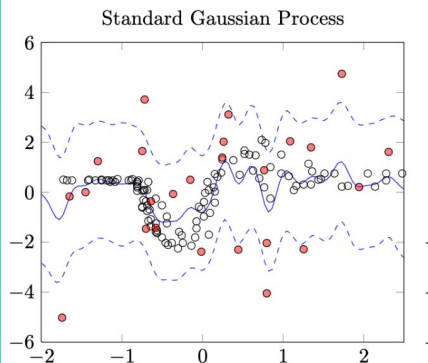


$$\begin{aligned} \text{ELBO}(q(\mathbf{f})) &= \int q(\mathbf{f}) \log p(\mathbf{y}|\mathbf{f}) d\mathbf{f} - \mathcal{KL}(q(\mathbf{f}), p(\mathbf{f})) \\ &= \sum_{i=1}^N \int q(f_i) \log p(y_i|f_i) d\mathbf{f} - \mathcal{KL}(q(\mathbf{f}), p(\mathbf{f})) \end{aligned}$$

$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$

SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

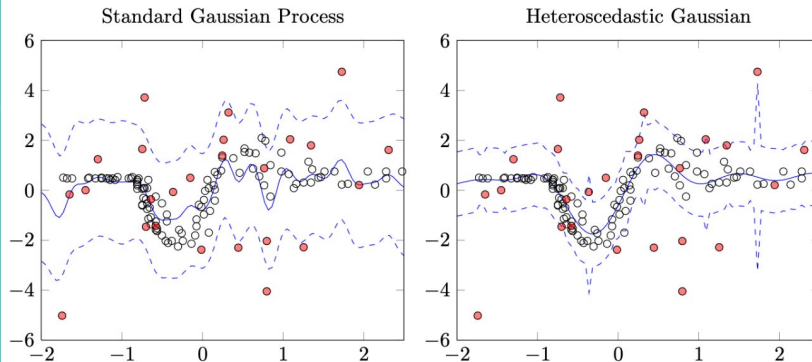


SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

~~$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$~~

$$y_i \sim \mathcal{N}(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)})$$



SVGPs for non-Gaussian data?

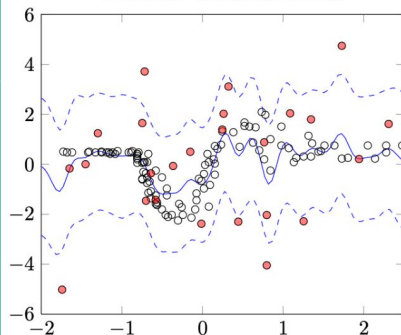
(Hensman et al. 2015, Saul et al. 2016)

~~$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$~~

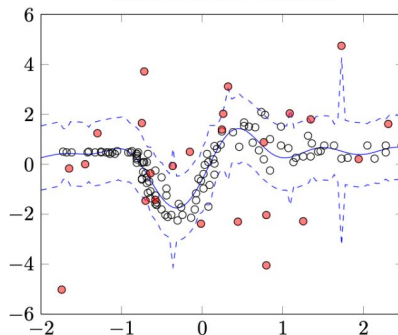
~~$$y_i \sim \mathcal{N}(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)})$$~~

$$y_i \sim \text{St}(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)}, \nu)$$

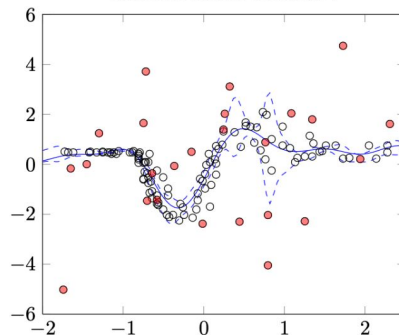
Standard Gaussian Process



Heteroscedastic Gaussian



Heteroscedastic Student-t



SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

ELBO($q(\mathbf{f}_0), q(\mathbf{f}_1)$)

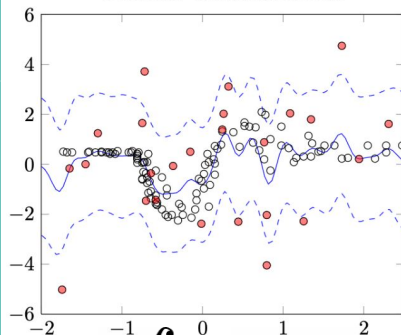
$$= \int q(\mathbf{f}_0)q(\mathbf{f}_1) \log p(\mathbf{y}|\mathbf{f}_0, \mathbf{f}_1) d\mathbf{f}_0 d\mathbf{f}_1 - \mathcal{KL}(q(\mathbf{f}_0), p(\mathbf{f})) - \mathcal{KL}(q(\mathbf{f}_1), p(\mathbf{f}))$$

~~$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$~~

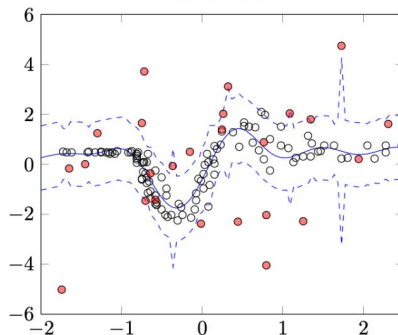
~~$$y_i \sim \mathcal{N}(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)})$$~~

$$y_i \sim St(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)}, \nu)$$

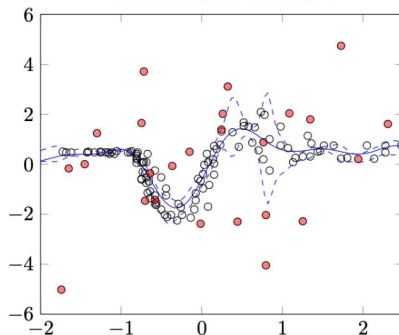
Standard Gaussian Process



Heteroscedastic Gaussian



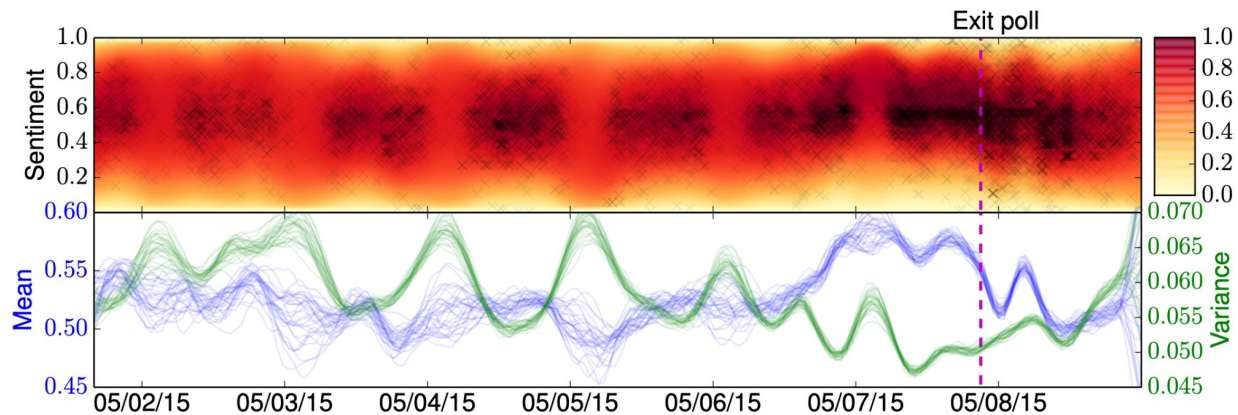
Heteroscedastic Student-t



SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

$$y_i \sim \mathcal{B}(\alpha = f_0(\mathbf{x}_i), \beta = e^{f_1(\mathbf{x}_i)})$$



GPs for
high-dim
data?



Beware the curse of
dimensionality



GPs for high-dim data?

- GPs are great in high-dim
- RBF kernels are not.....
- $l_i \propto \sqrt{D}$

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2l^2}}$$

GPs for
high-dim
data?

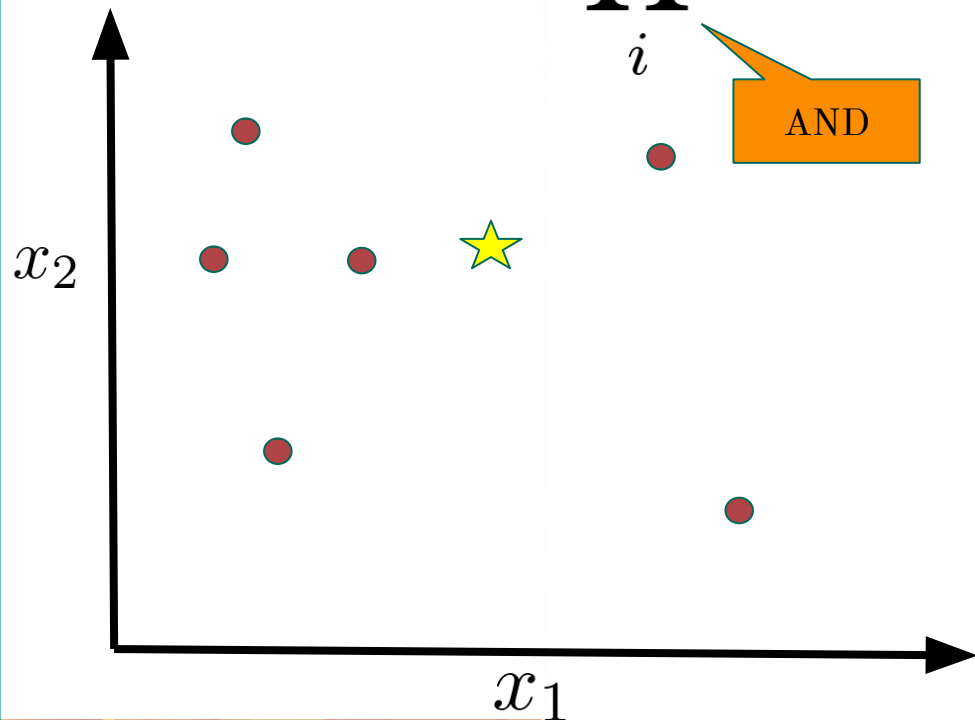
$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$



AND

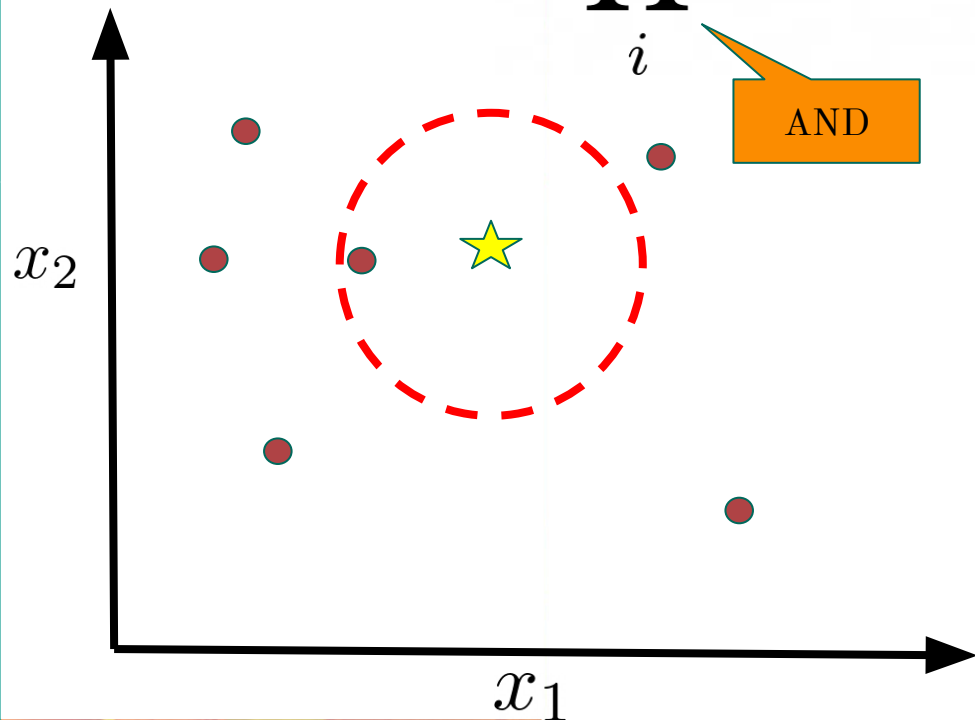
GPs for high-dim data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$

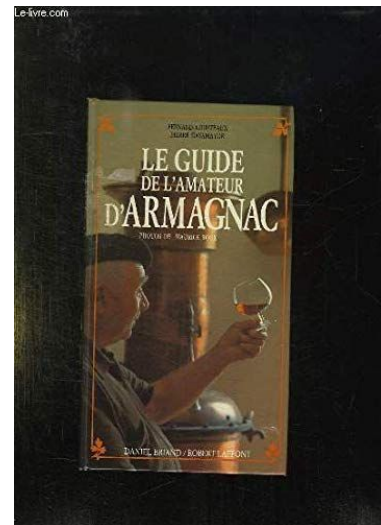


GPs for high-dim data?

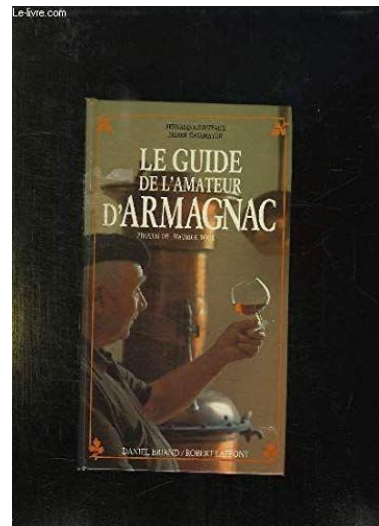
$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$



GPs for high-dim data?

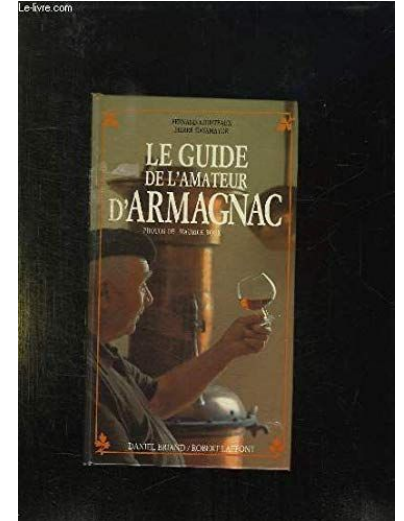


GPs for high-dim data?



GPs for high-dim data?

- Type of still (column/pot?)
- Type of grape (Ugni Blanc?)
- Wood for the barrel
- Location (Armagnac-Ténarèze, Bas-Armagnac ,Haut-Armagnac?)
- Blend
- Age



GPs for
high-dim
data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$



AND

GPs for high-dim data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$

AND

$$k_1(\mathbf{x}, \mathbf{y}) = \sum_i^d k_i(x_i, y_i)$$

OR

$$k_2(\mathbf{x}, \mathbf{y}) = \sum_{i < j}^d k_i(x_i, y_i) k_j(x_j, y_j)$$

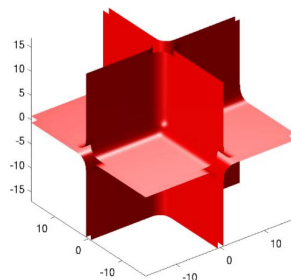
Additive Gaussian Processes

$$k(x, y) = k_0 + \sum k_i(x_i, y_i) + \sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)$$

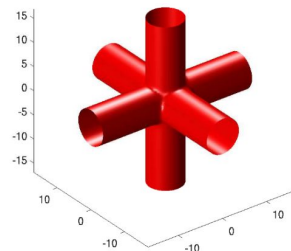
Additive Gaussian Processes

$$k(x, y) = \overset{0}{k_0} + \sum \overset{1}{k_i(x_i, y_i)} + \sum_{i < j} \overset{2}{k_i(x_i, y_i)k_j(x_j, y_j)}$$

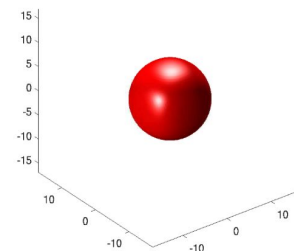
(Duvenaud et al 2011)



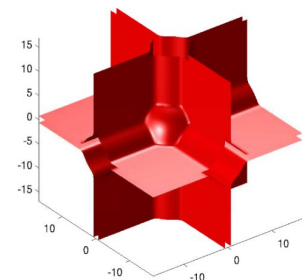
1st order interactions
 $k_1 + k_2 + k_3$



2nd order interactions
 $k_1k_2 + k_2k_3 + k_1k_3$



3rd order interactions
 $k_1k_2k_3$
(Squared-exp kernel)



All interactions
(Additive kernel)

Additive Gaussian Processes

$$k(x, y) = k_0 + \sum k_i(x_i, y_i) + \sum_{i < j} k_i(x_i, y_i)k_j(x_j, y_j)$$



Ginsbourger et al. (2016)

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

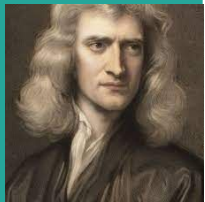
Additive Gaussian Processes

$$k(x, y) = \overset{0}{k_0} + \sum \overset{1}{k_i(x_i, y_i)} + \sum_{i < j} \overset{2}{k_i(x_i, y_i)k_j(x_j, y_j)}$$

\Updownarrow

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

- Standard RBF $\rightarrow O(d(N^2 + NM))$
- d additive RBF $\rightarrow O(2^d(N^2 + NM))$



Additive Gaussian Processes

- Newton Girard (Duvenaud et al 2011)

$$k(x, y) = \overset{0}{k_0} + \sum \overset{1}{k_i(x_i, y_i)} + \sum_{i < j} \overset{2}{k_i(x_i, y_i)k_j(x_j, y_j)}$$
$$\Updownarrow$$
$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

- Standard RBF $\rightarrow O(d(N^2 + NM))$
- d additive RBF $\rightarrow O(2^d(N^2 + NM))$
- d additive BBF (NG) $\rightarrow O(d^2(N^2 + NM))$

Additive Gaussian Processes

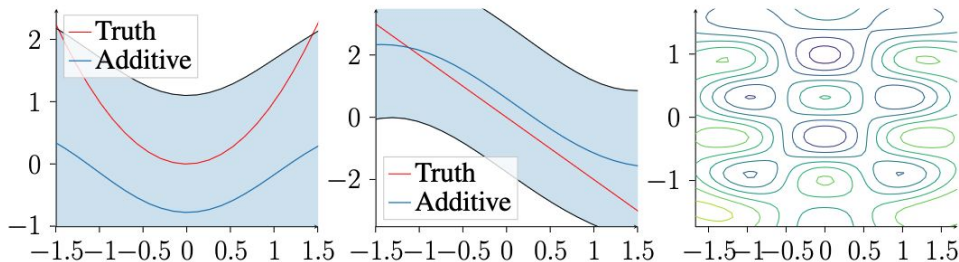
$$k(x, y) = k_0 + \sum k_i(x_i, y_i) + \sum_{i < j} k_i(x_i, y_i)k_j(x_j, y_j)$$



Ginsbourger et al. (2016)

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



$$E[f_i(x_i) | \mathcal{D}] = k_i(x_i, X) K(X, X)^{-1} \mathbf{y}$$

Additive Gaussian Processes

- Orthogonalise (Durrande et al 2012)

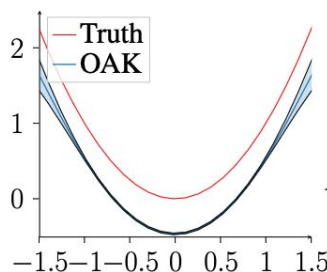
$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

$$k(x, y) = \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)}$$

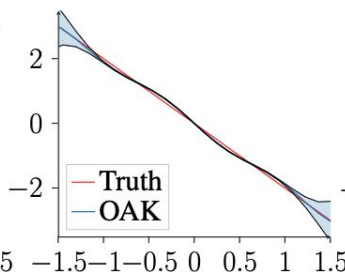


$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

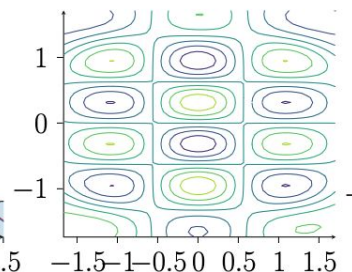
$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



(f) f_1



(g) f_2



(h) Interaction

Additive Gaussian Processes

- Orthogonalise (Durrande et al 2012)

$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

- By conditioning

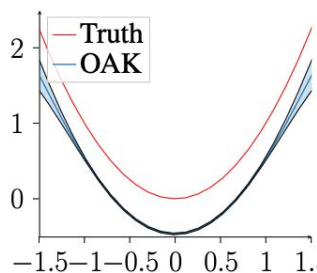
$$f_i(x_i) \Big| \int f_i(x_i) p(x_i) dx_i = 0$$

$$k(x, y) = \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)}$$

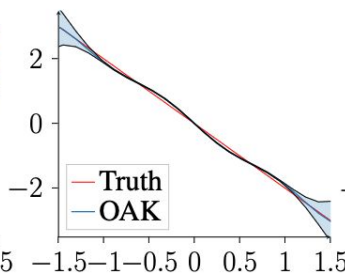
$$\Updownarrow$$

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

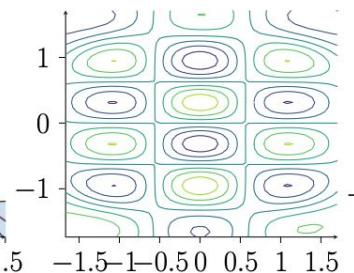
$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



(f) f_1



(g) f_2



(h) Interaction

Additive Gaussian Processes

- Orthogonalise (Durrande et al 2012)

$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

- By conditioning

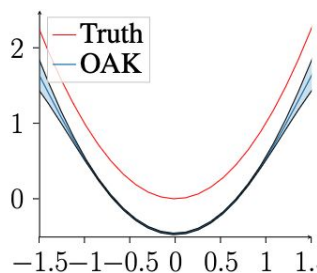
$$f_i(x_i) \Big| \int f_i(x_i) p(x_i) dx_i = 0$$

$$k(x, y) = \sum_{i=1}^n \alpha_i k_i(x_i, y_i) + \sum_{i < j} \alpha_{ij} k_{ij}(x_i, x_j, y_i, y_j)$$

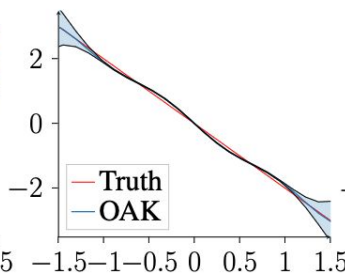
This model is quite interpretable.....

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

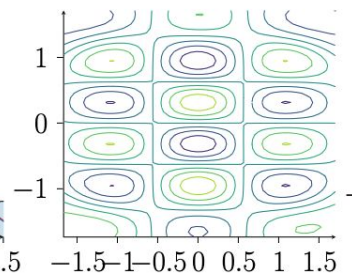
$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



(f) f_1



(g) f_2



(h) Interaction

So, lets learn an equation

—

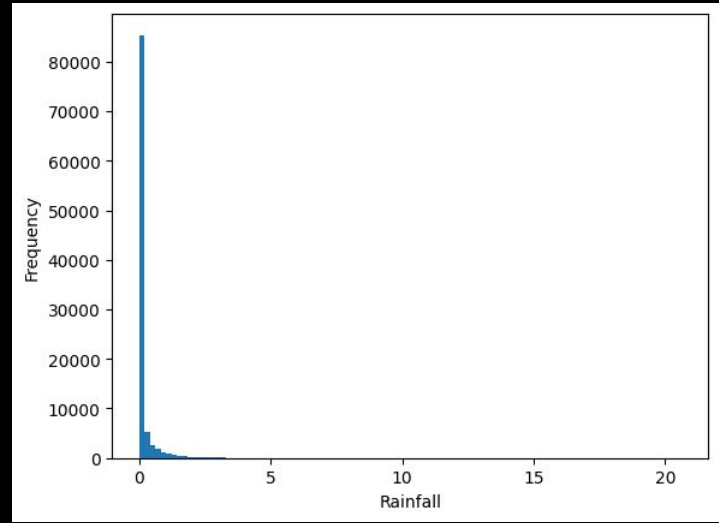
Predicting rainfall

- >100 climate variables → rainfall



Predicting rainfall

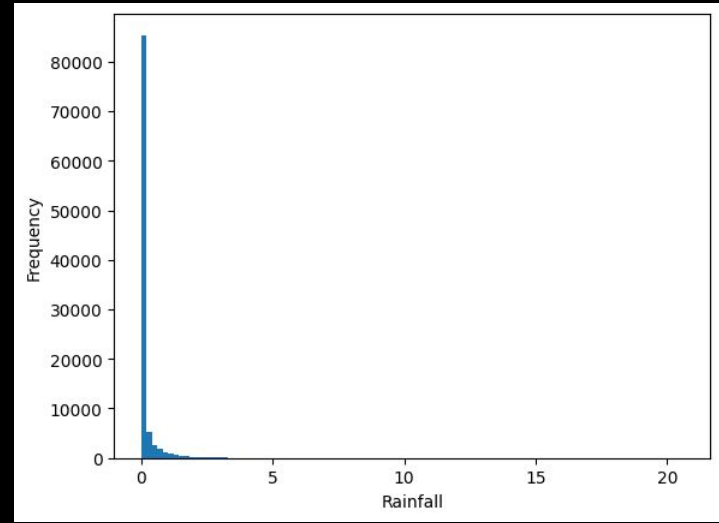
- >100 climate variables → rainfall
- Non-Gaussian (Bernoulli-gamma)



$$\cancel{p(y|f) = \mathcal{N}(f, \sigma^2)}$$

Predicting rainfall

- >100 climate variables → rainfall
- Non-Gaussian (Bernoulli-gamma)



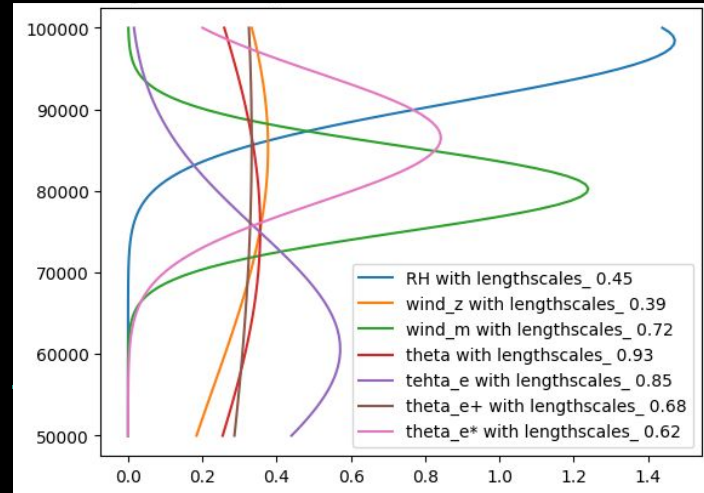
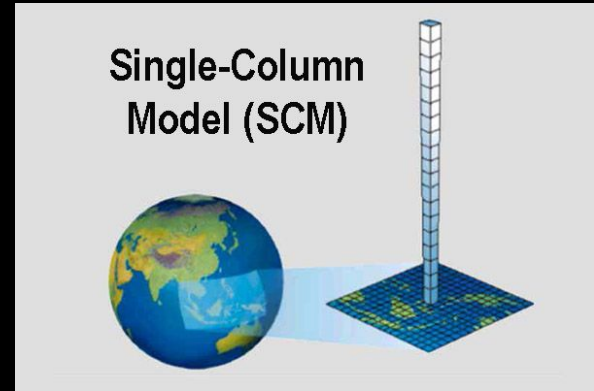
~~$$p(y|f) = \mathcal{N}(f, \sigma^2)$$~~

$$p(\underline{y|f_1, f_2, f_3}) = \begin{cases} 1 - f_1 & \text{if } y = 0 \\ f_1 \Gamma(f_2, f_3) & \text{o.w.} \end{cases}$$

Predicting rainfall

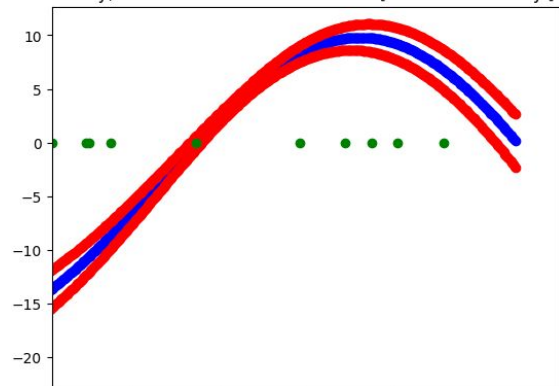
- >100 climate variables → rainfall
- Non-Gaussian (Bernoulli-gamma)
- Data-driven vertical integration

<https://e3sm.org/single-column-model-intercomparison-of-diurnal-cycle-of-precipitation/>



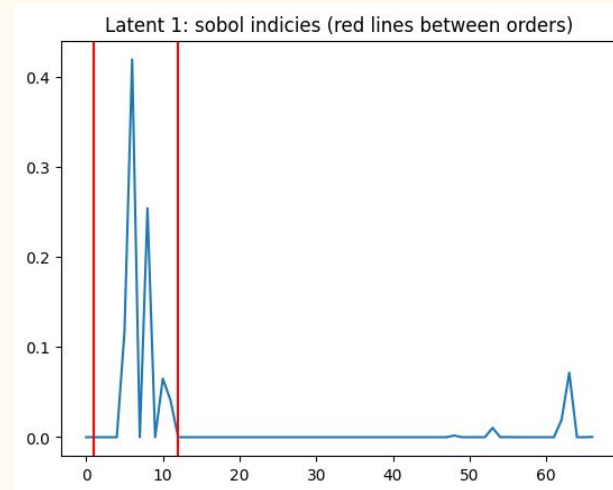
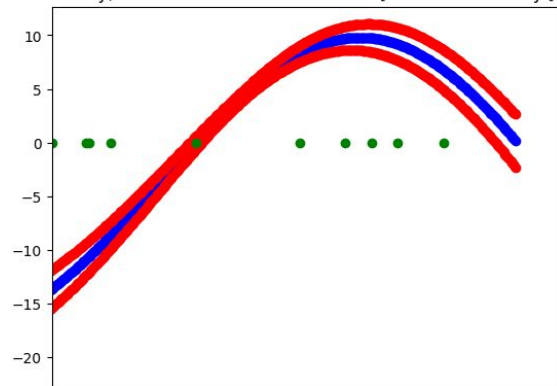
Additive GP model output

Latent 0 rank 0: Best guess (and uncertainty) at additive contributions from ['Relative Humidity']with sobol index 0.581364255678434



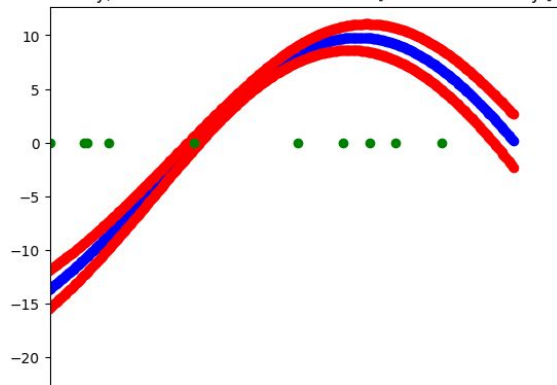
Additive GP model output

Latent 0 rank 0: Best guess (and uncertainty) at additive contributions from ['Relative Humidity']with sobol index 0.581364255678434

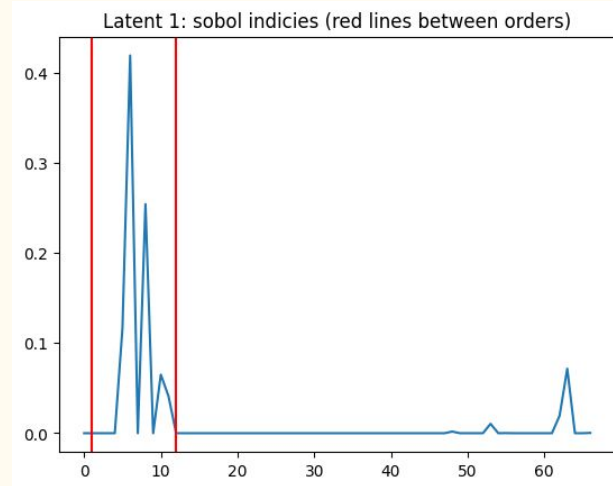
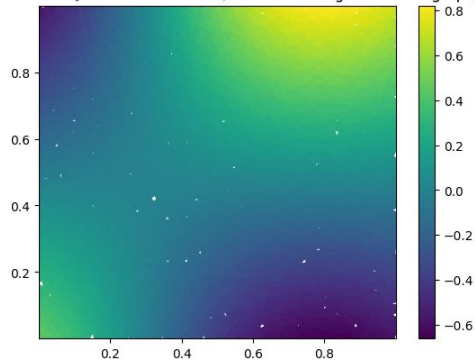


Additive GP model output

Latent 0 rank 0: Best guess (and uncertainty) at additive contributions from ['Relative Humidity'] with sobol index 0.581364255678434



Latent 1 rank 1: Best guess at additive contribution from ['Sensible heat flux', 'Stdev of sub-gridscale orography'] with sobol index 0.0715865896060103



Learn a Stochastic Eq (via lots of easy SRs)

$$p(y|f_1, f_2, f_3) = \begin{cases} 1 - f_1 & \text{if } y = 0 \\ f_1 \Gamma(f_2, f_3) & \text{o.w.} \end{cases}$$



Just
illustrative
results

- $f_1 = e^{\lambda_0 + \lambda_1 RH - \lambda_2 RH \sigma_0}$
- $f_2 = \lambda_3 + \lambda_4 (SHF - \lambda_5)^2$
- $f_3 = \lambda_6 + \lambda_7 \theta_+$

Learn a Stochastic Eq (via lots of easy SRs)

$$p(y|f_1, f_2, f_3) = \begin{cases} 1 - f_1 & \text{if } y = 0 \\ f_1 \Gamma(f_2, f_3) & \text{o.w.} \end{cases}$$



Just
illustrative
results

$$E[y] = e^{\lambda_0 + \lambda_1 RH - \lambda_2 RH * \sigma_0} \frac{(SHF - \lambda_3)^2}{\lambda_4 + \lambda_5 \theta_+}$$

Whats next

- Gravity waves / cloud cover
- **Learn the likelihood** structure (another layer of symbolic regression)
- Improve “orthogonality” for correlated inputs
- **Sample** multiple candidate equations (Pareto front?)
- More user interaction
- Encode **known physics** (symmetries, invariances, conservation laws e.t.c)

La presqu'île de Blackpool



Thanks for your
time!

Extra slides

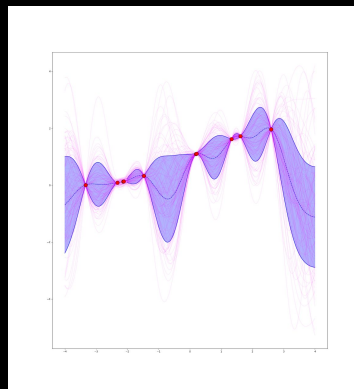


Scientific priors via conditioning

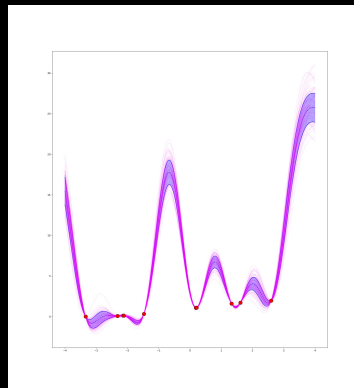
Condition on an integral

$$O(f) = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

$\mathbb{P}(f)$



$\mathbb{P}(f|O)$

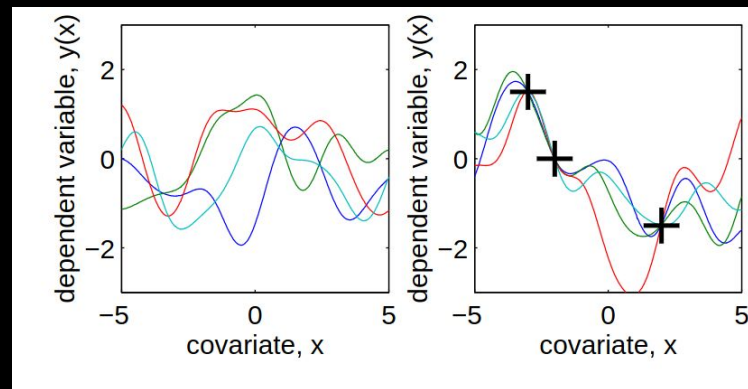


Scientific priors via conditioning

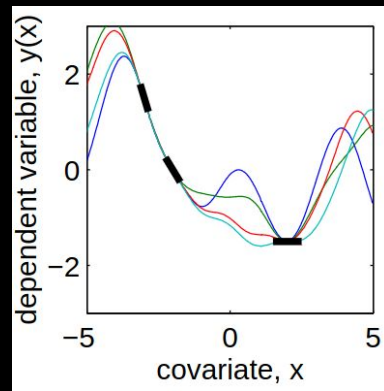
Condition on an derivative

$$O(\mathbf{f}) = \frac{\partial f}{\partial \mathbf{x}}$$

$\mathbb{P}(f)$



$\mathbb{P}(f|O)$



Scientific priors via conditioning

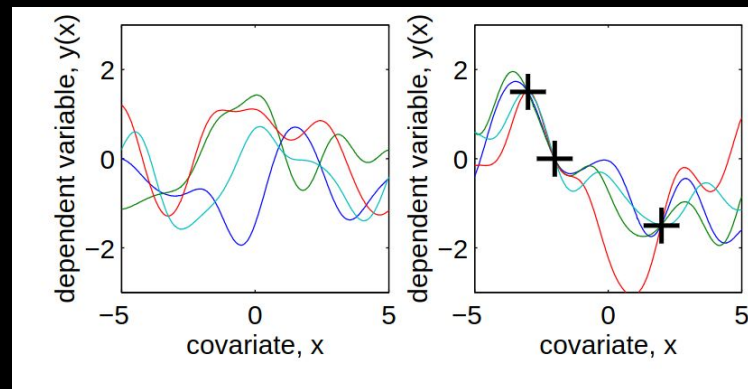
Condition on an derivative



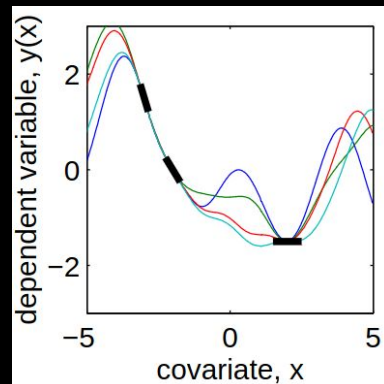
Padidar et al. (2021)

$$O(\mathbf{f}) = \frac{\partial f}{\partial \mathbf{x}}$$

$\mathbb{P}(f)$



$\mathbb{P}(f|O)$



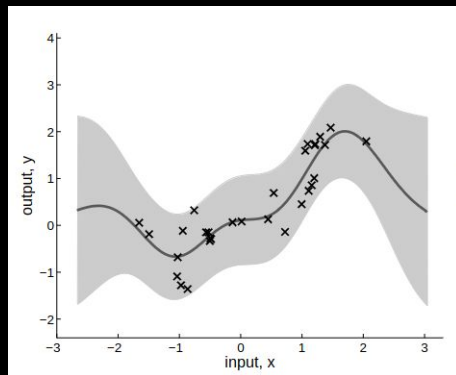
Solak et al. (2002)

Scientific priors via conditioning

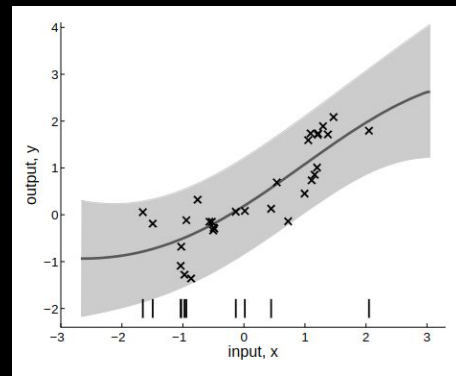
Condition on monotonicity

$$O(\mathbf{f}) = \left(\frac{\partial f}{\partial \mathbf{x}} > 0 \right)$$

$$\mathbb{P}(f)$$



$$\mathbb{P}(f|O)$$

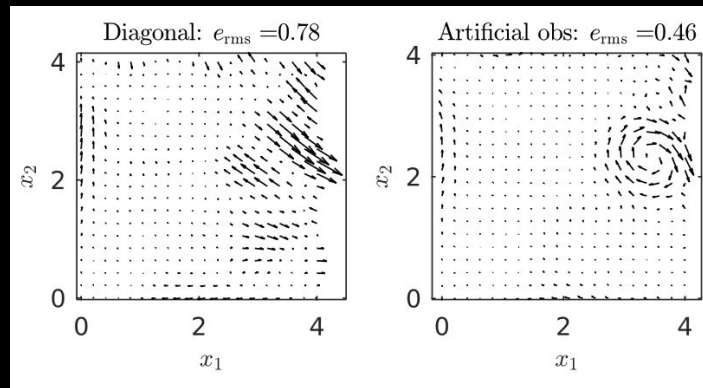


Scientific priors via conditioning

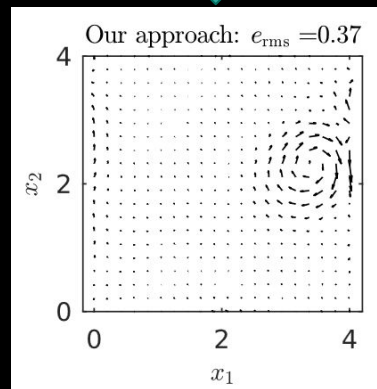
Condition on linear operator

$$O(f) = \frac{\partial f}{\partial x_1} + \frac{\partial f}{\partial x_2}$$

$\mathbb{P}(f)$



$\mathbb{P}(f|O)$

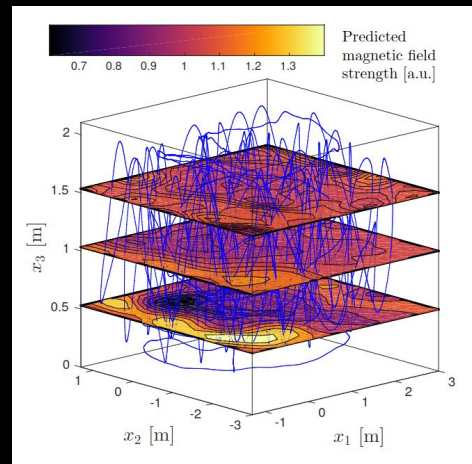


Scientific priors via conditioning

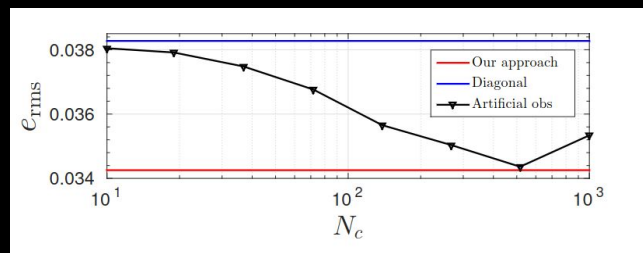
Condition on linear operator

$$\mathcal{O}(\mathbf{f}) = \nabla \times \mathbf{f}$$

$$\mathbb{P}(\mathbf{f})$$



$$\mathbb{P}(\mathbf{f} | \mathcal{O})$$

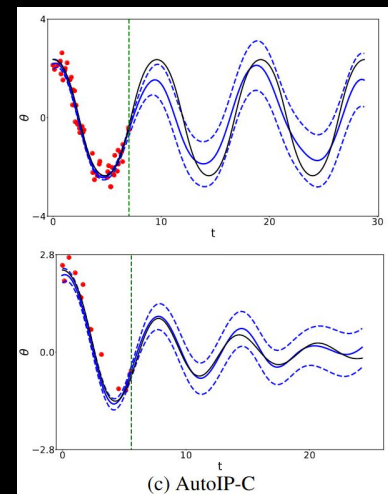
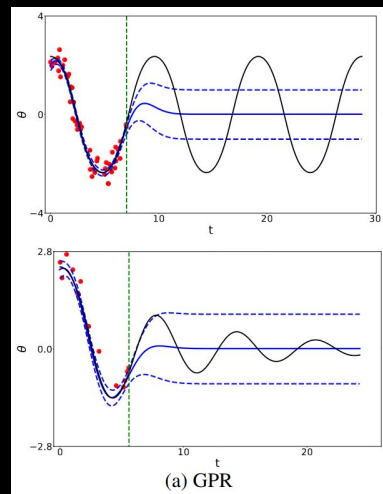


Scientific priors via conditioning

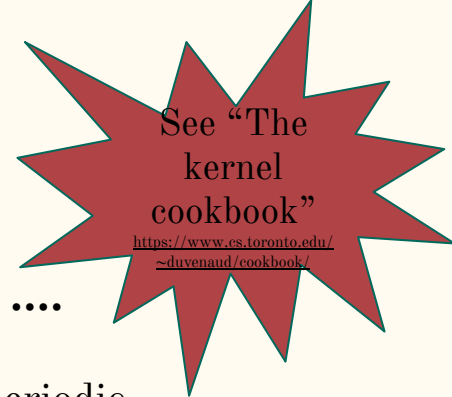
Condition on whatever you want and
pretend its Gaussian

$$\mathbb{P}(f|D) \propto \mathbb{P}(D|f)\mathbb{P}(f)\mathbb{P}(O(f))$$

$$\frac{df^2}{dt} + \sin(t) + \beta \frac{df}{dt} = 0$$

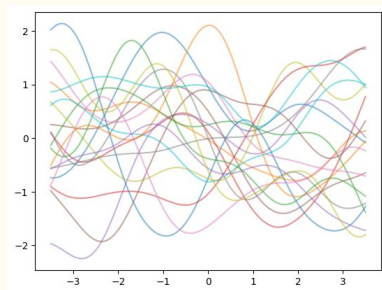


Two ways to be encode info into GPs

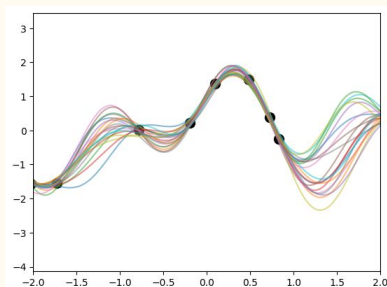


1) Additional conditioning

$\mathbb{P}(f)$



$\mathbb{P}(f|\mathcal{D})$

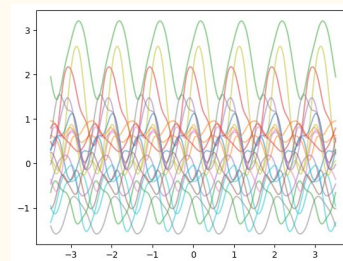


2) Thinking hard

- I want f to be periodic
- So choose a periodic kernel



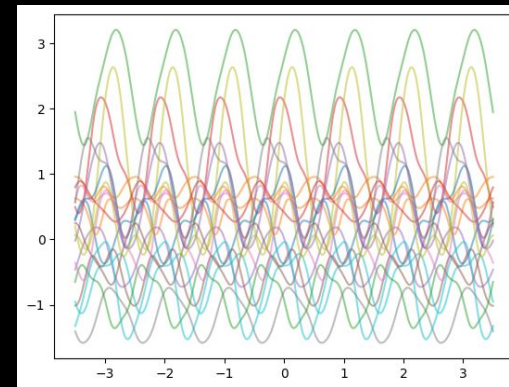
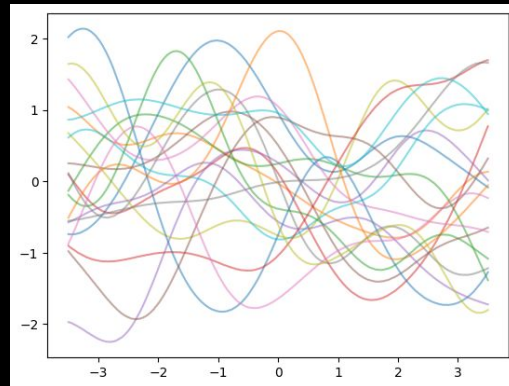
$$k_{per}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(\frac{-2 \sin^2(\pi|\mathbf{x} - \mathbf{x}'|/p)}{l^2}\right)$$



Scientific priors via kernel design

Fiddle with the kernel to get
periodicity

$$k_{per}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(\frac{-2 \sin^2(\pi|\mathbf{x} - \mathbf{x}'|/p)}{l^2}\right)$$



Scientific priors via kernel design

General idea:

$$T(f) = f \Leftrightarrow T(k(\mathbf{x}, \cdot)) = k(\mathbf{x}, \cdot)$$

$$\hat{k}(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') + k(T(\mathbf{x}), \mathbf{x}')$$

Ginsbourger et al. 2013

Van der Wilk et al. 2018

