

Gaussian Processes & Regressors

A visual introduction to Gaussian Processes

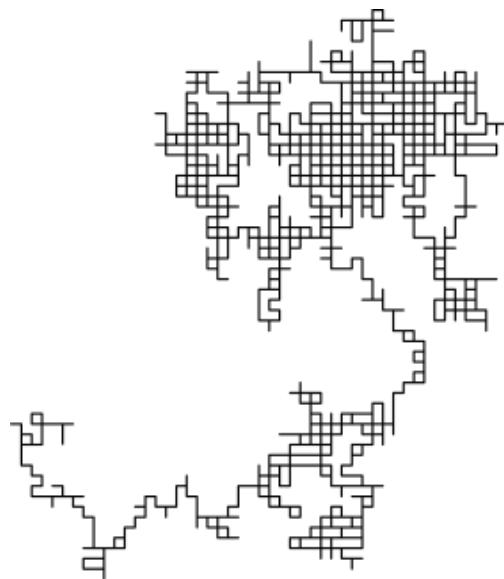
Mads-Peter V. Christiansen

Stochastic Process:

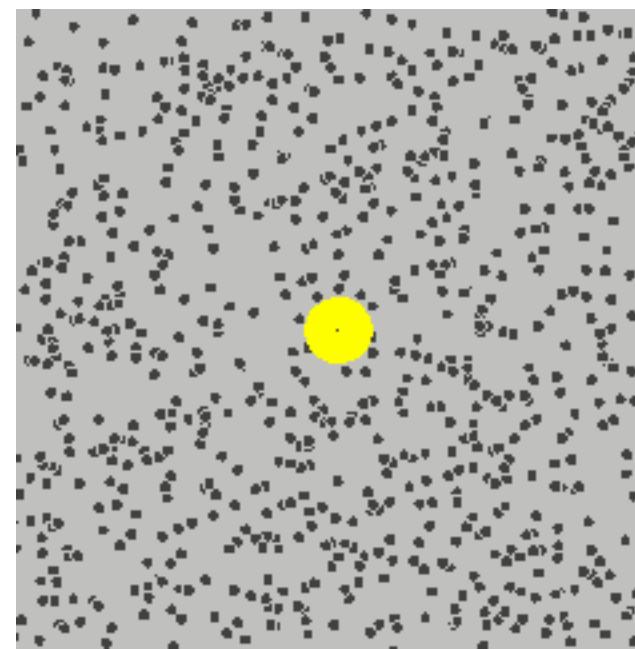
- "A collection of random variables indexed by time, space, or some other domain."

Think a **recipe** for generating sequences of random variables.

Random walk



Brownian Motion

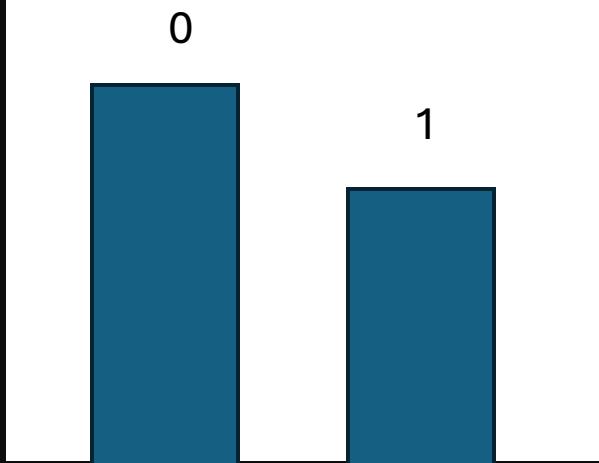


Bernoulli Process:

- Sequence of random variables each drawn iid. from a Bernoulli distribution.

This definition tells us what the **recipe** is.

Bernoulli Distribution



Tells you how to draw a single random variable

Bernoulli Process

Sample paths drawn from the process

$$X = \{0,0,1,0,1,1,0 \dots \}$$

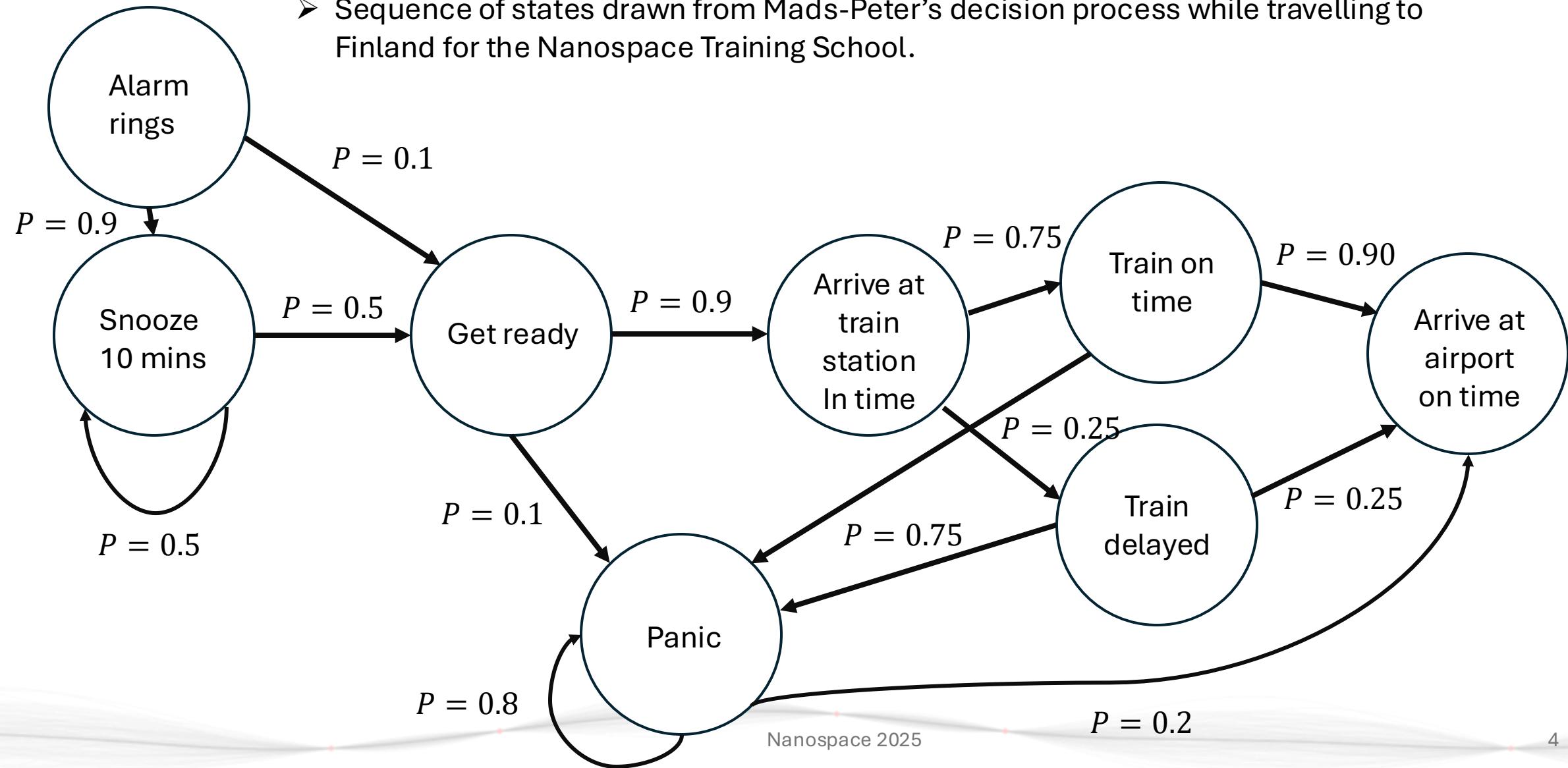


$$X = \{1,1,1,0,0,0,1 \dots \}$$

Tells you how to draw a collection of random variables

Mads-Peter Process

- Sequence of states drawn from Mads-Peter's decision process while travelling to Finland for the Nanospace Training School.



Mads-Peter Process

- Sequence of states drawn from Mads-Peter's decision process while travelling to Finland for the Nanospace Training School.

Samples from my decision process



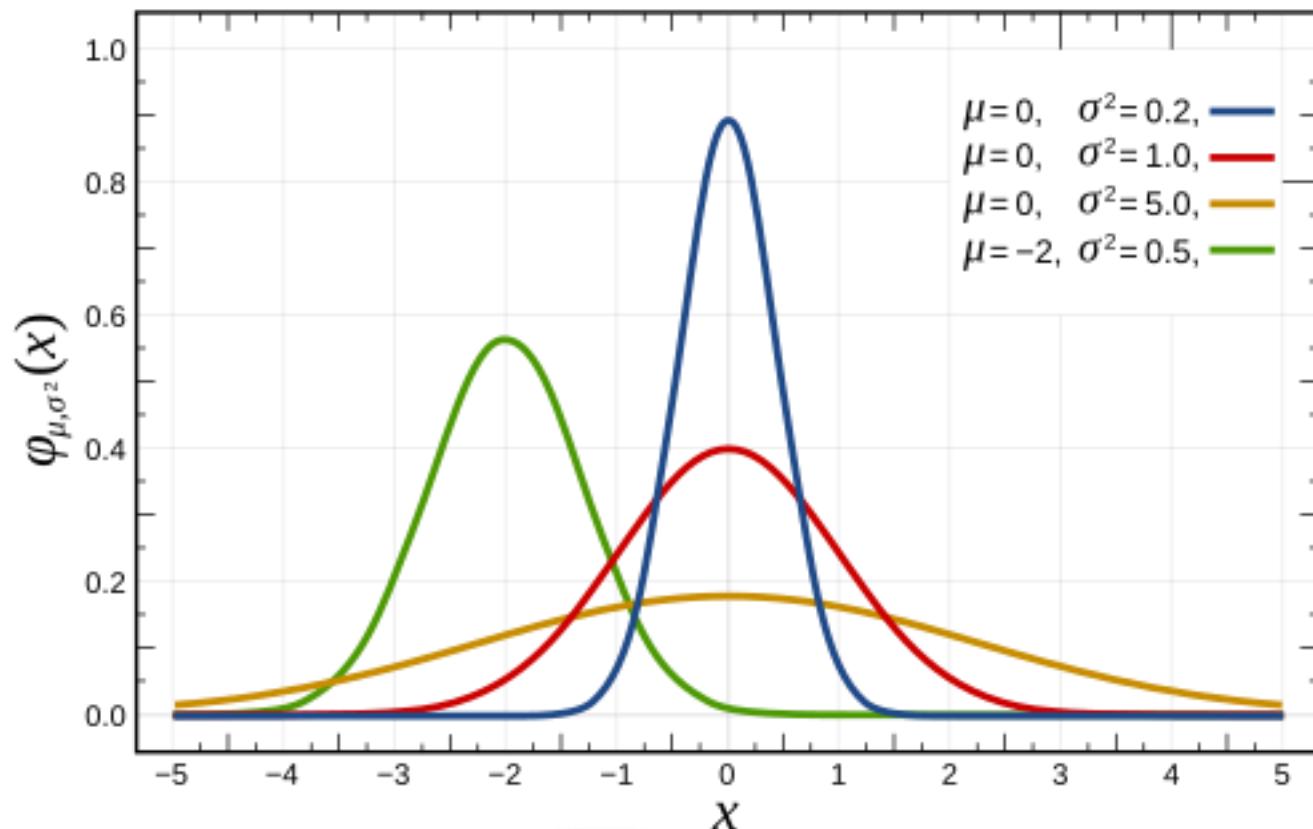
Gaussian

- **Carl Friedrich Gauss(ian)**: German mathematician, astronomer, geodesist and physicist.
- **Gauss(ian) (unit)**: Unit of magnetic flux
- **Gaussian function**: A function of the form e^{-x^2}
- **Gaussian distribution**: Continuous probability distribution.



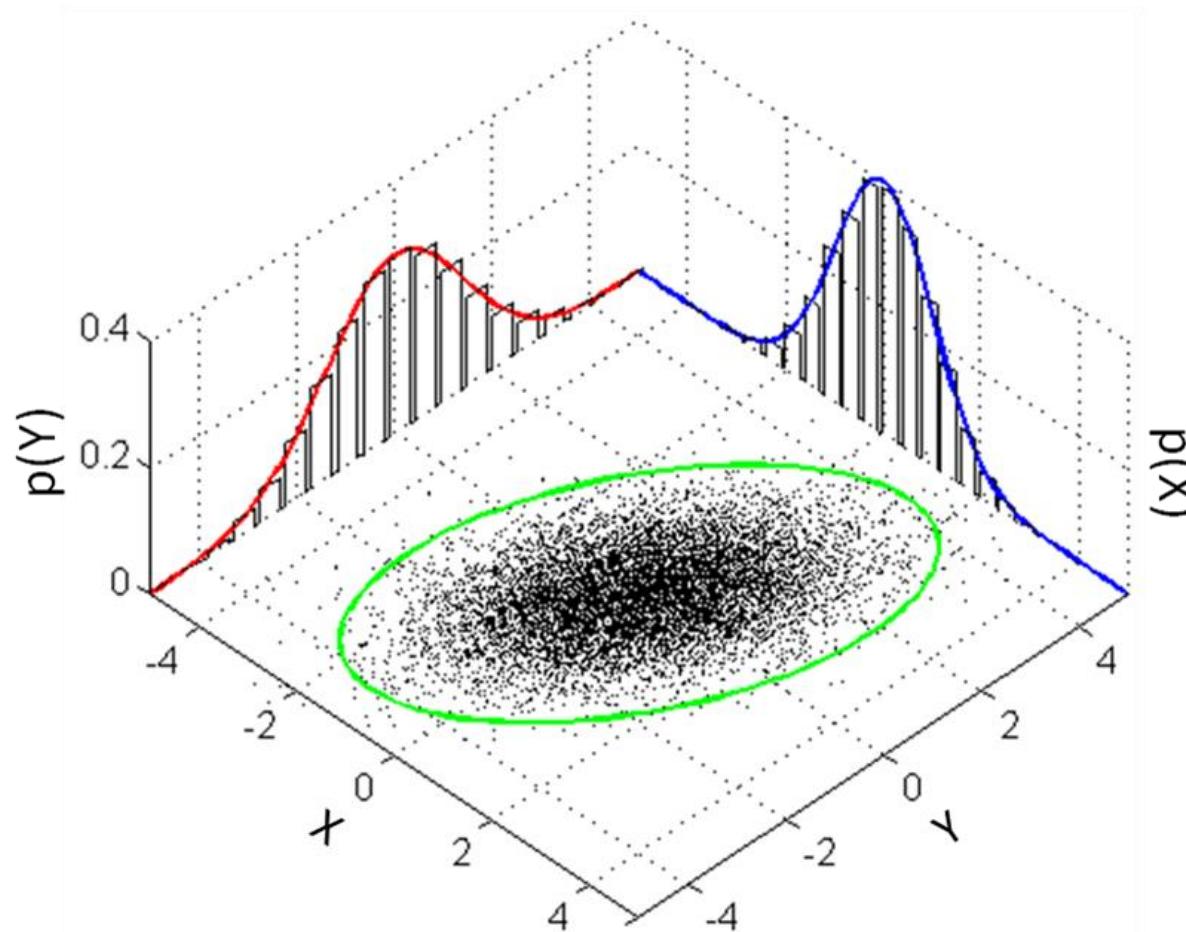
Gaussian distribution

- Normal distribution with the familiar bell shape parameterized by the mean μ and the variance σ^2 .



Multivariate Gaussian distribution

- A distribution over vectors, parameterized by the mean μ and the covariance Σ .

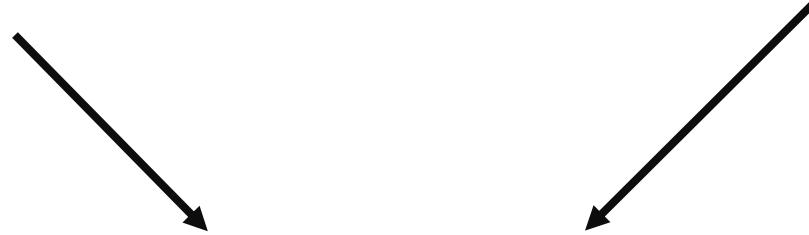


Multivariate Gaussian distribution

- Multivariate normal distribution parameterized by the mean μ and the covariance Σ .

Stochastic Process:

- "A collection of random variables indexed by time, space, or some other domain."



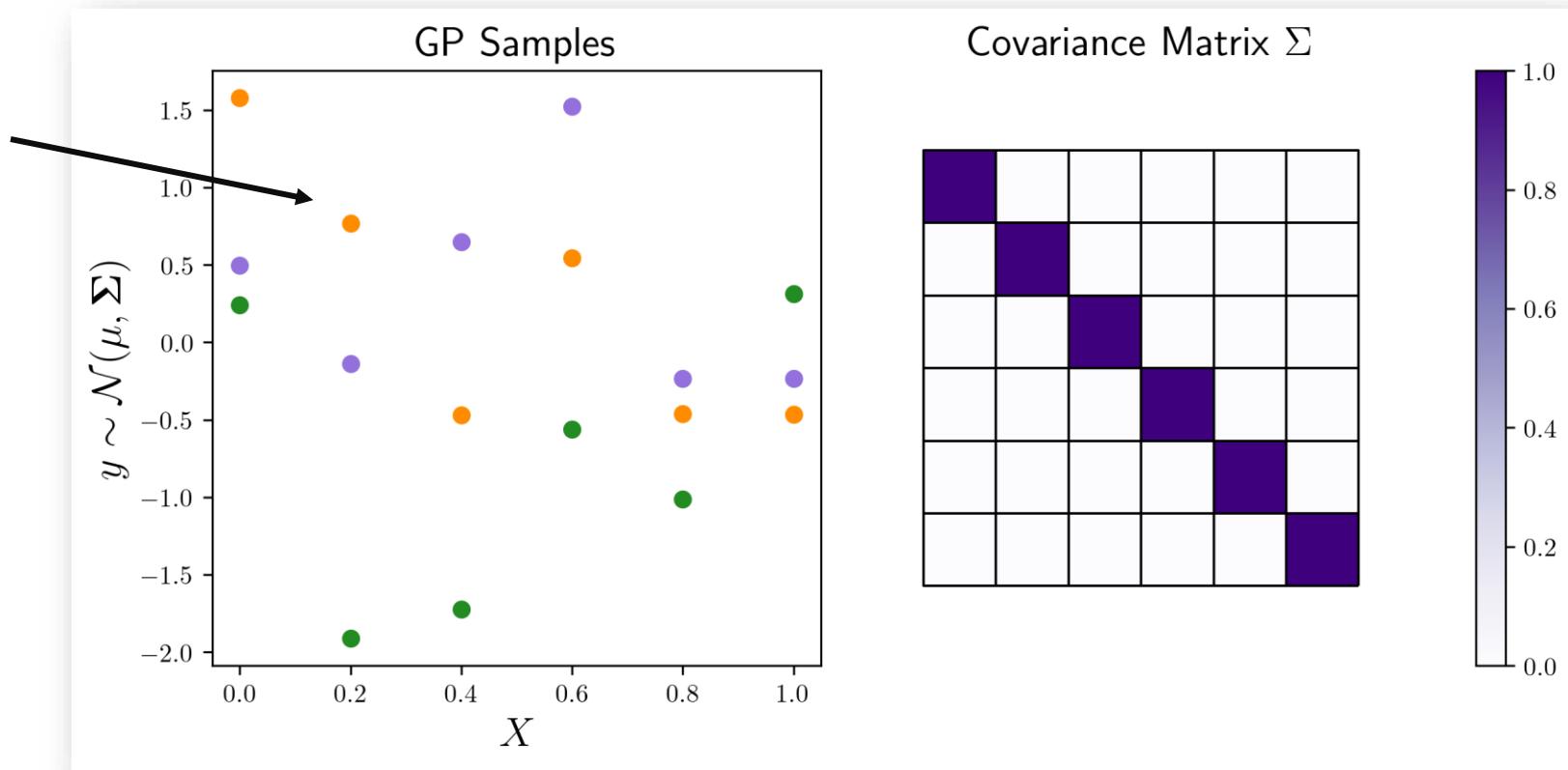
Gaussian Process

- A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.

Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”

Three samples from the GP.

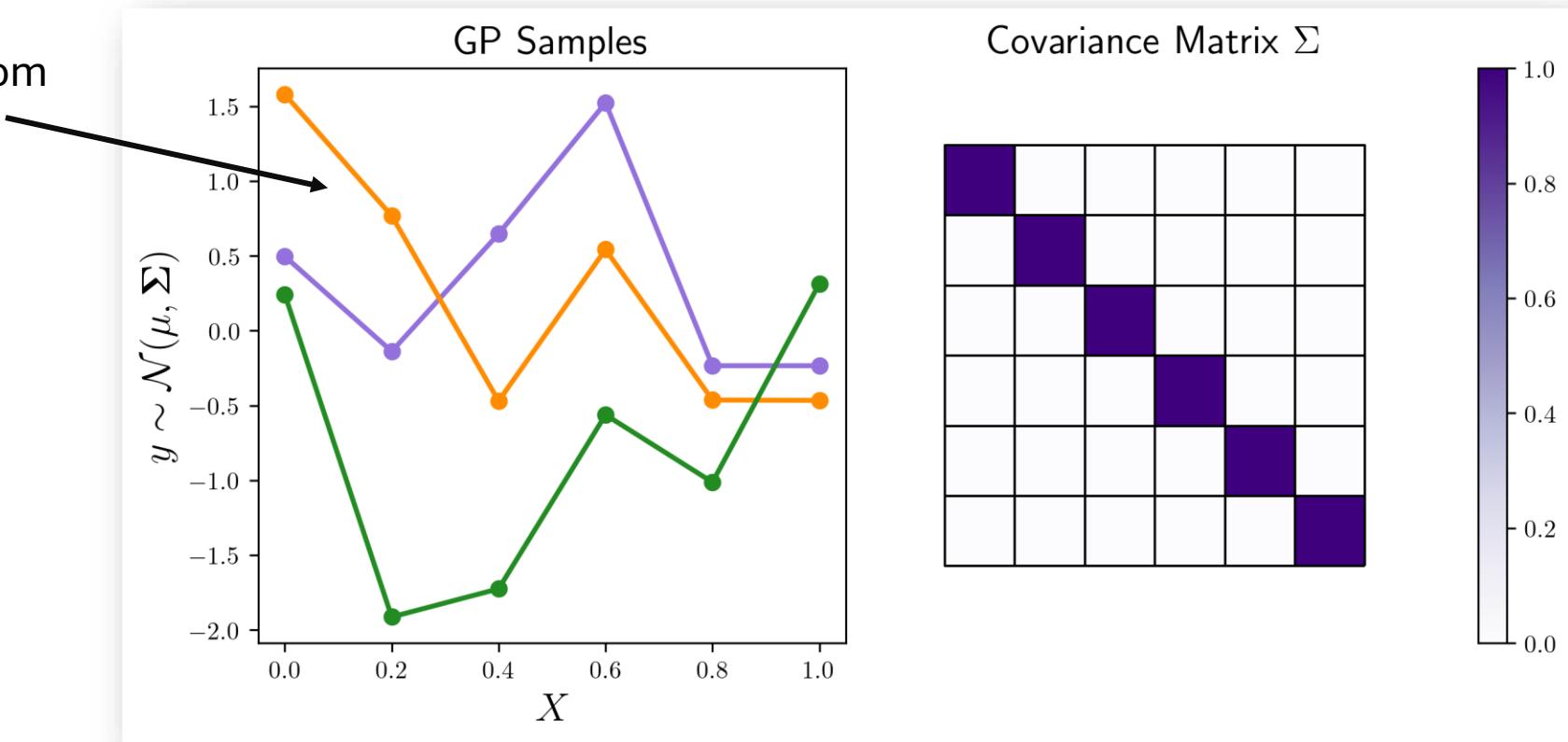


“... indexed by time, space or some other domain.”

Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”

Three samples from the GP.



“... indexed by time, space or some other domain.”

Gaussian Process

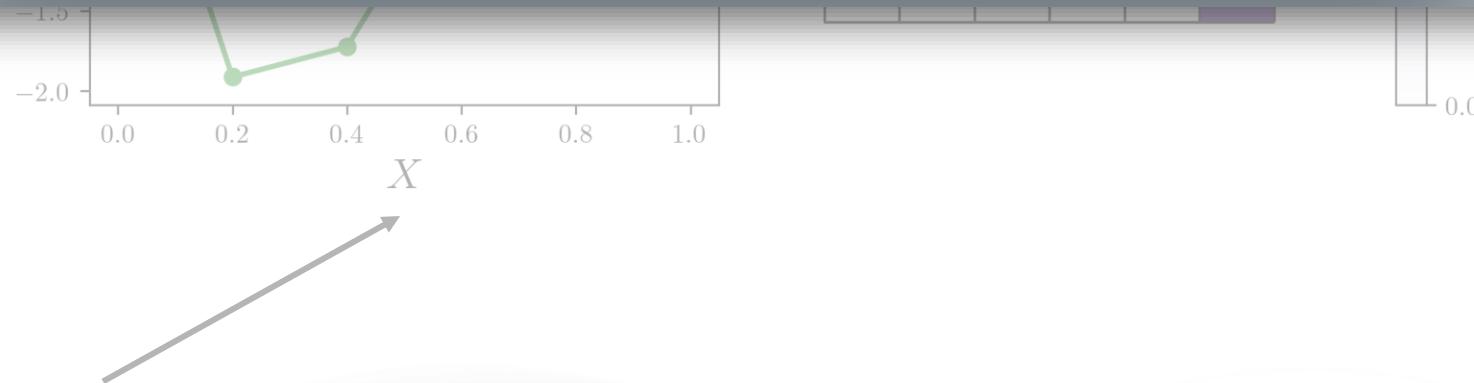
- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”

Three samples from the GP.

GP Samples

Covariance Matrix Σ

```
1 import numpy as np
2
3 X = np.array([0.0, 0.2, 0.4, 0.6, 0.8, 1.0]) # Sample points
4 covariance = np.eye(len(X))                  # Construct covariance matrix
5 mean = np.zeros(len(X))                     # Mean vector for the GP
6
7 sample = np.random.multivariate_normal(mean, covariance) # Sample from the GP
```

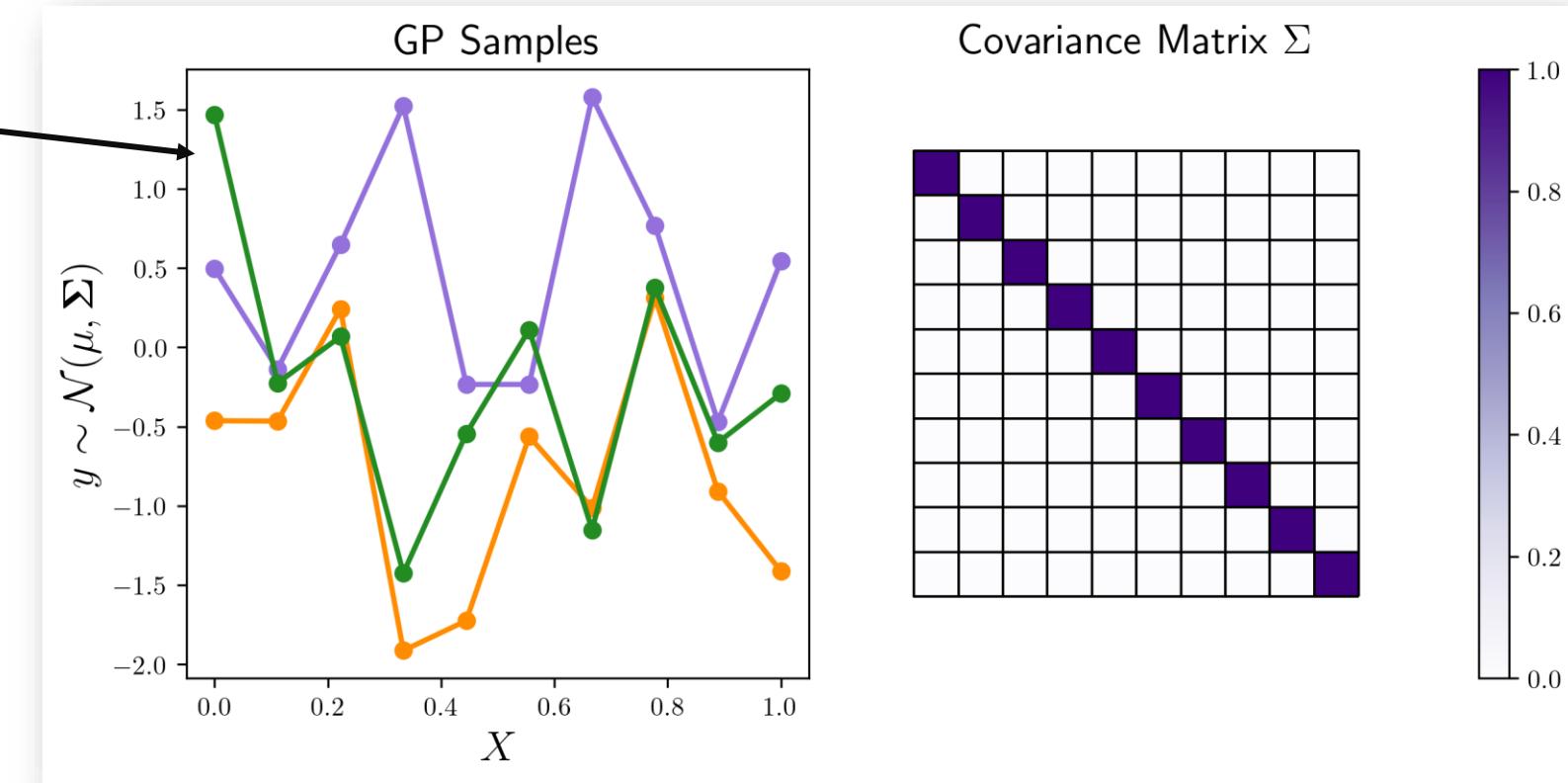


“... indexed by time, space or some other domain.”

Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”

Three samples from the GP.

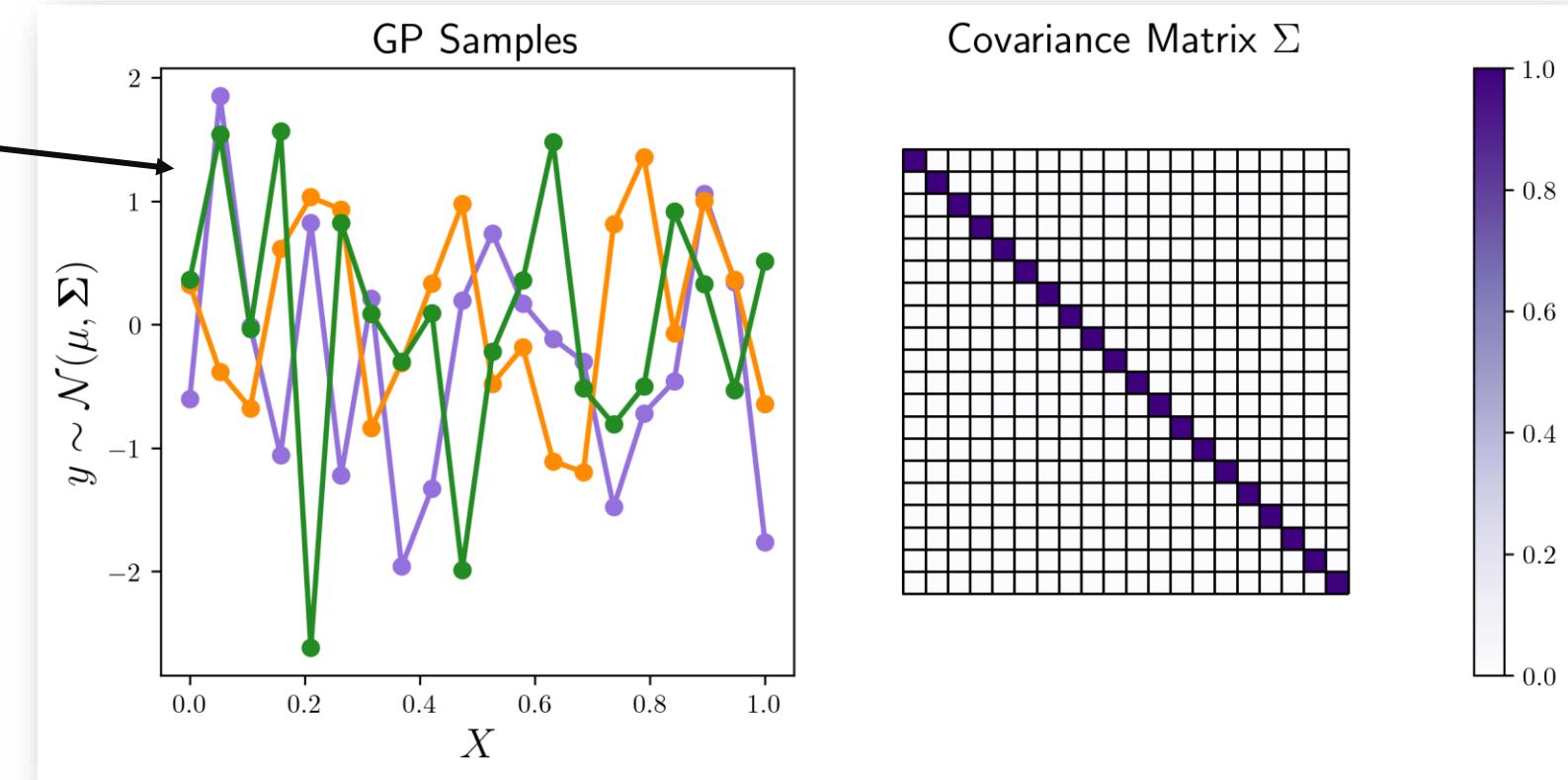


“... indexed by time, space or some other domain.”

Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”

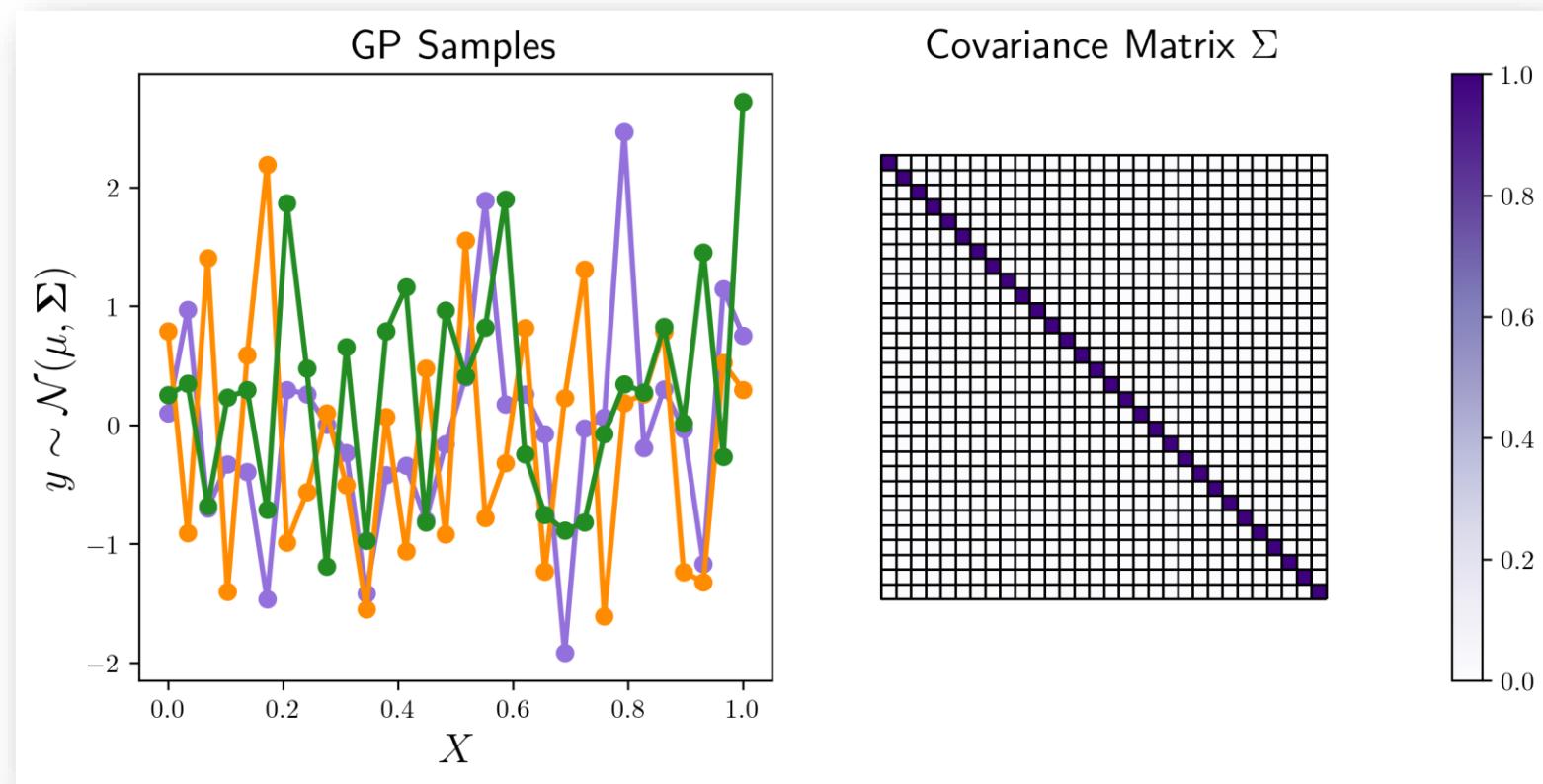
Three samples from the GP.



“... indexed by time, space or some other domain.”

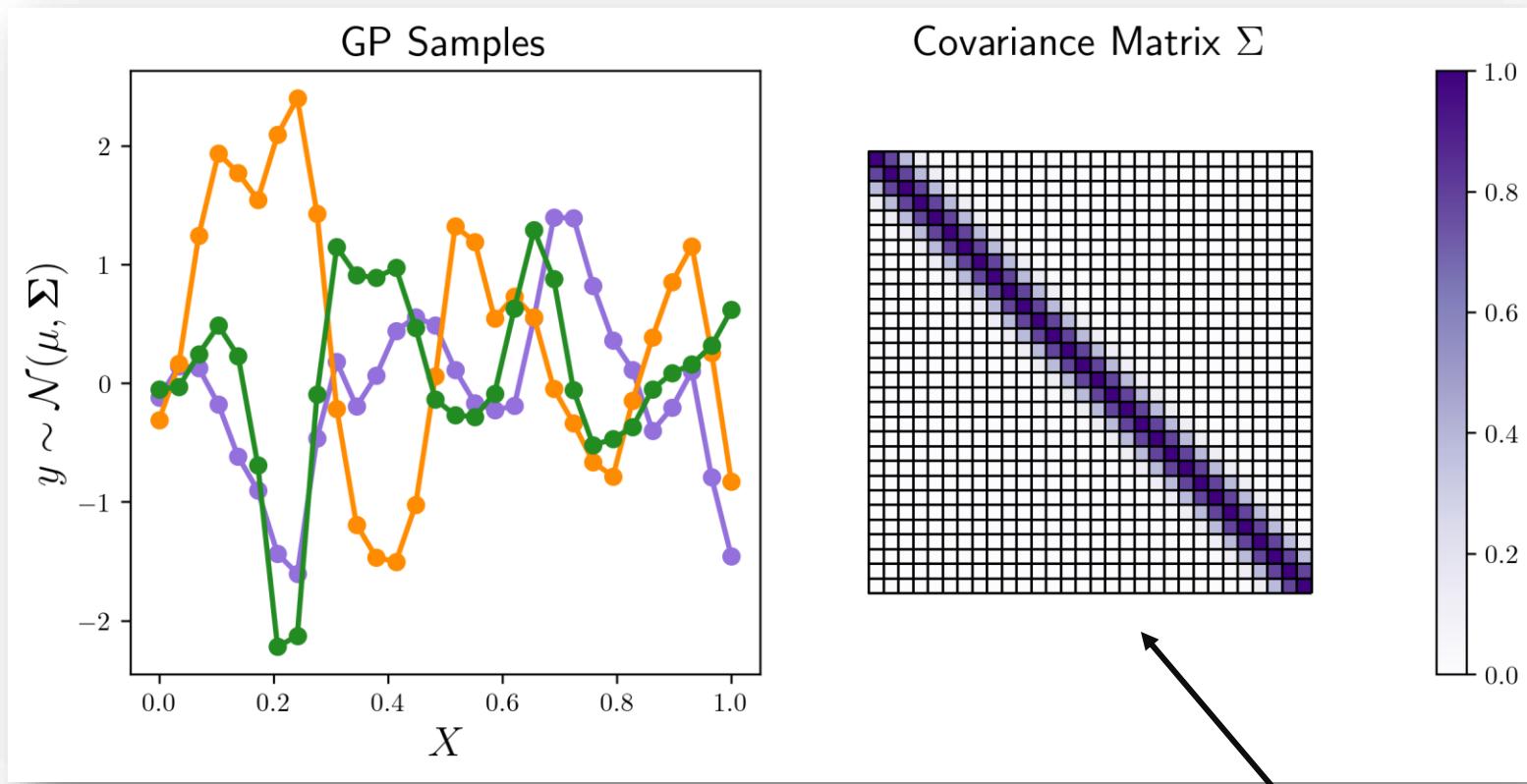
Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”



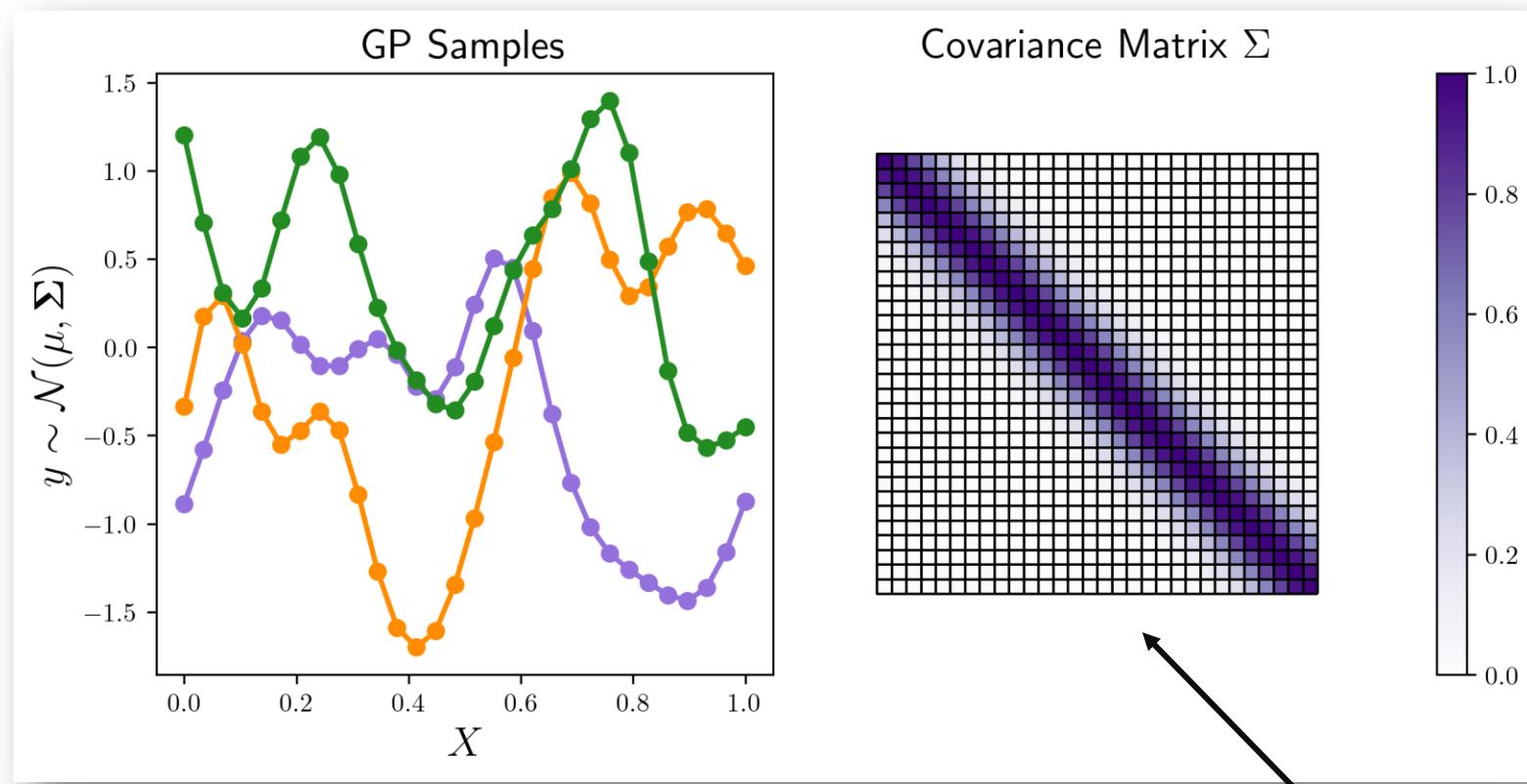
Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”



Gaussian Process

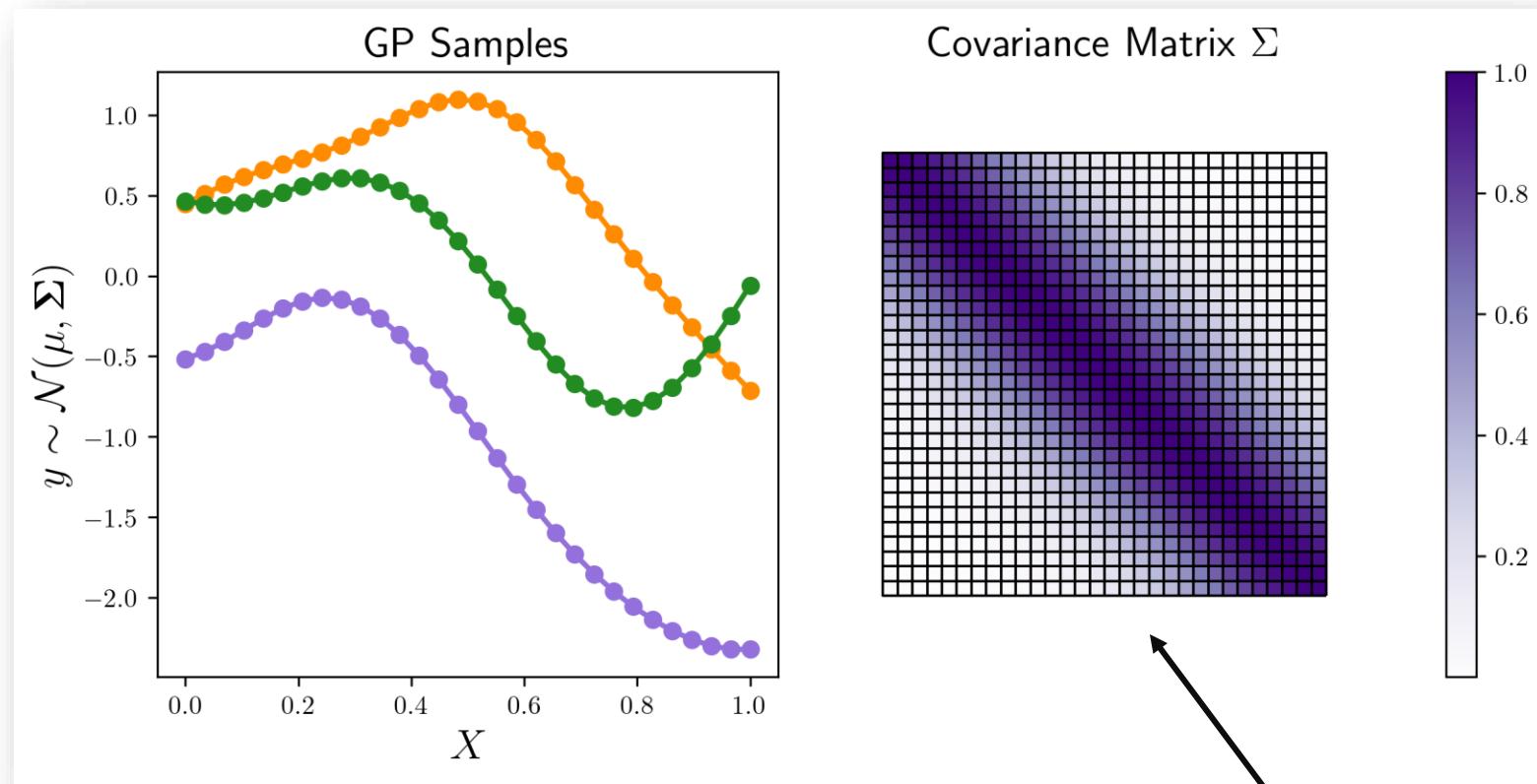
- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”



More structure in the covariance matrix

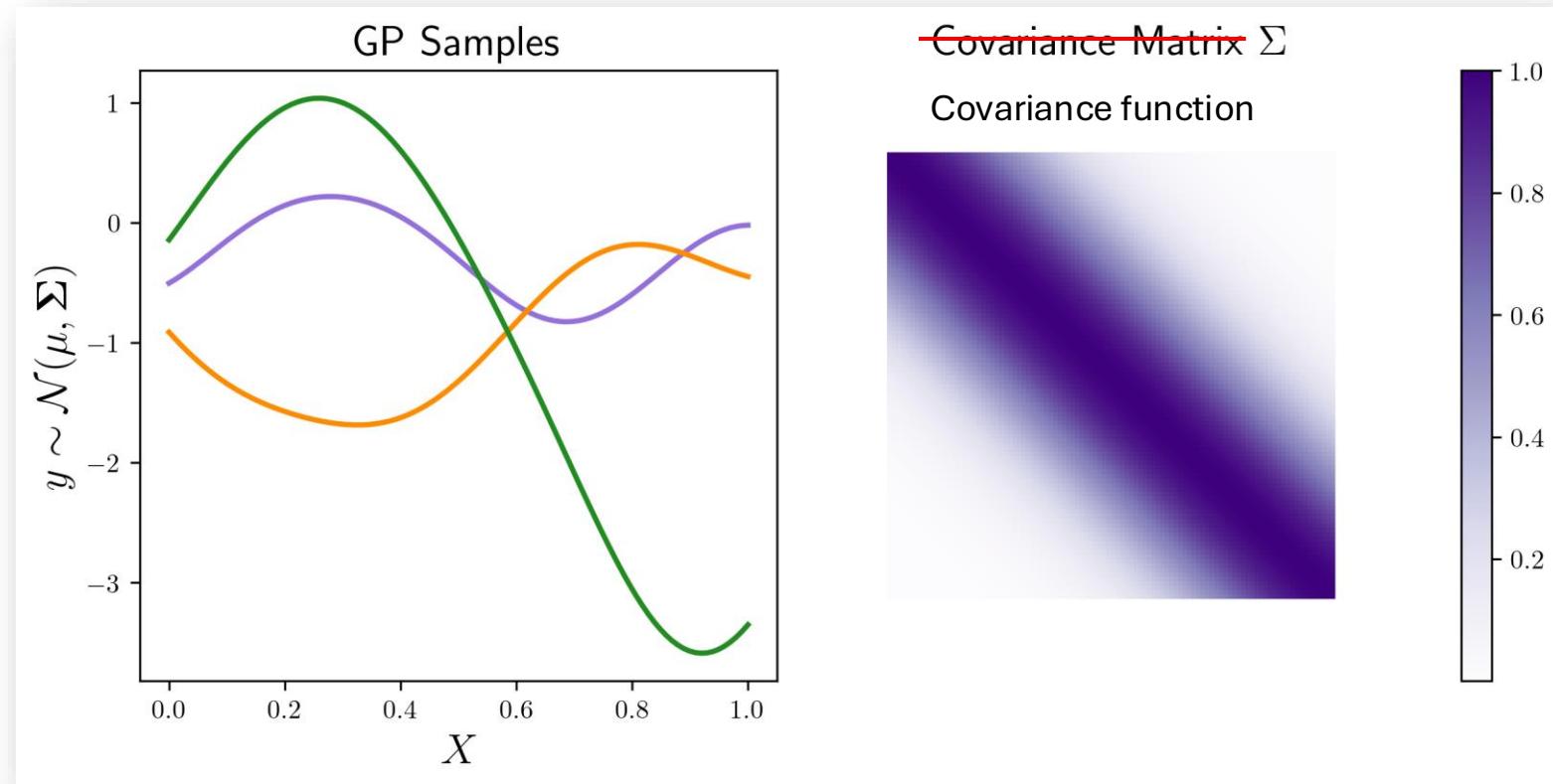
Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”



Gaussian Process

- “A stochastic process where samples are generated by drawing from a multivariate Gaussian distribution indexed by time, space or some other domain.”



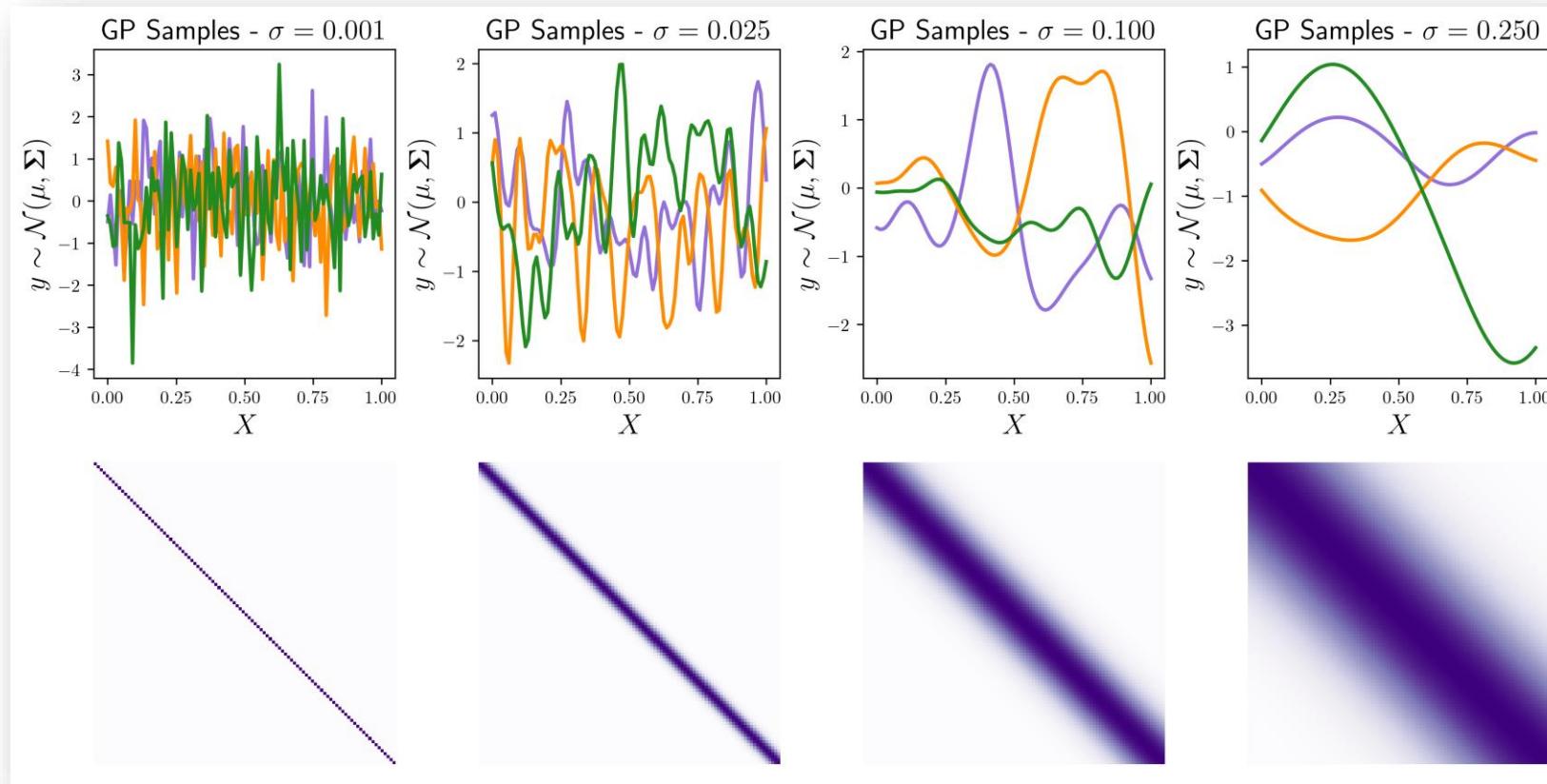
Gaussian Process

- “A probability distribution over *functions*, where any finite collection of function values is jointly *Gaussian distributed*.”

Covariance matrix → Covariance function = Kernel function

- “A kernel is a function that defines the covariance between function values at any pair of input locations.”

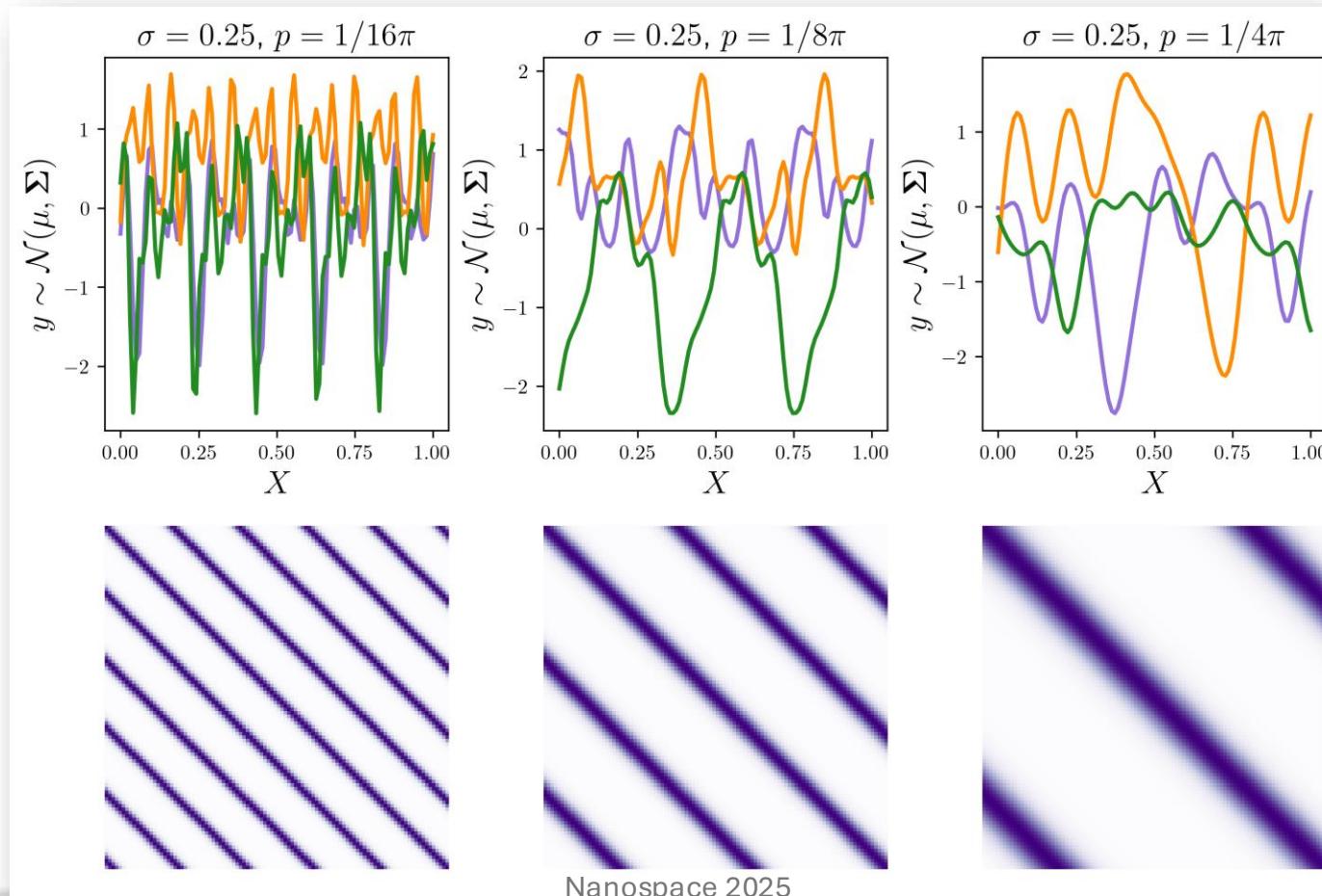
$$K(x_1, x_2) = \exp\left(-\frac{\|x_2 - x_1\|^2}{2\sigma^2}\right) \quad \text{RBF or Gaussian kernel.}$$



Covariance matrix → Covariance function = Kernel function

- “A kernel is a function that defines the covariance between function values at any pair of input locations.”

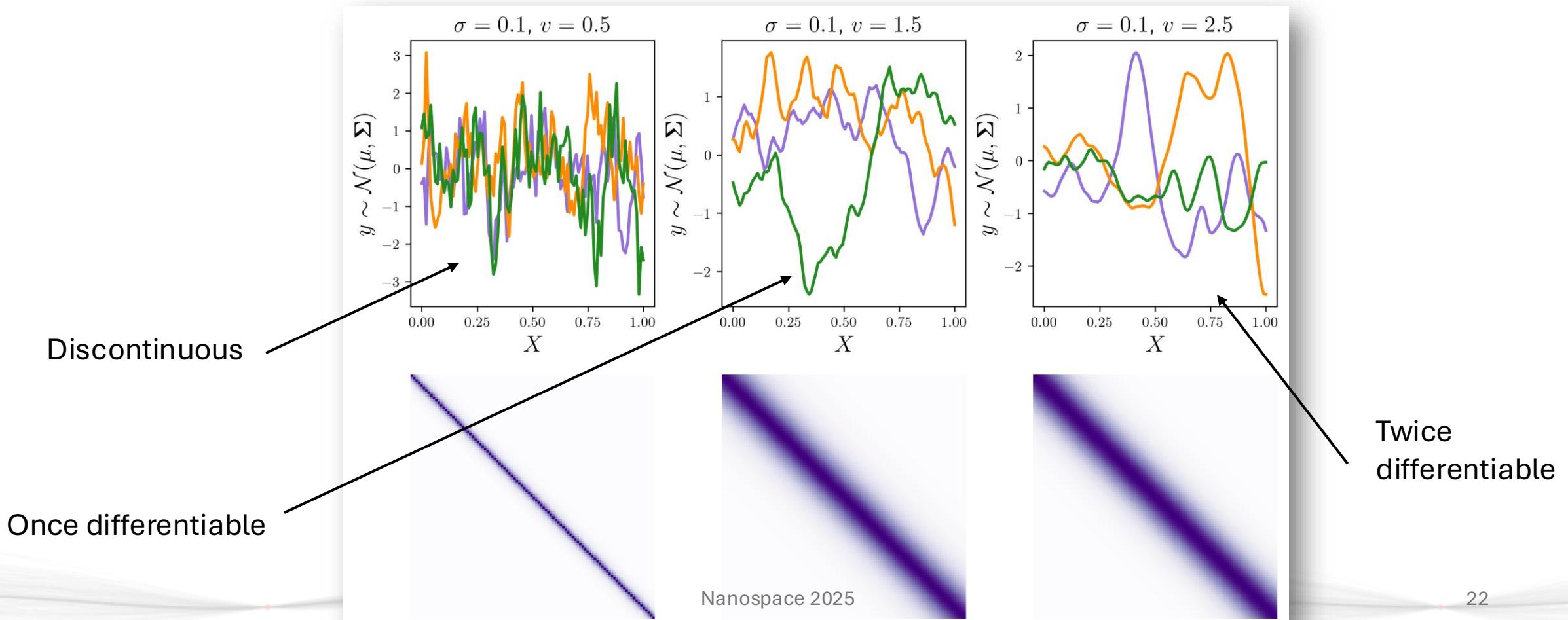
$$K(x_1, x_2) = \exp\left(-\frac{2}{\sigma^2} \sin^2\left(\frac{\pi||x_1 - x_2||}{p}\right)\right) \quad \text{Periodic kernel}$$



Covariance matrix → Covariance function = Kernel function

- “A kernel is a function that defines the covariance between function values at any pair of input locations.”

$$k(x_i, x_j) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{\sqrt{2\nu}}{l} d(x_i, x_j) \right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu}}{l} d(x_i, x_j) \right), \quad \text{Matern kernel}$$

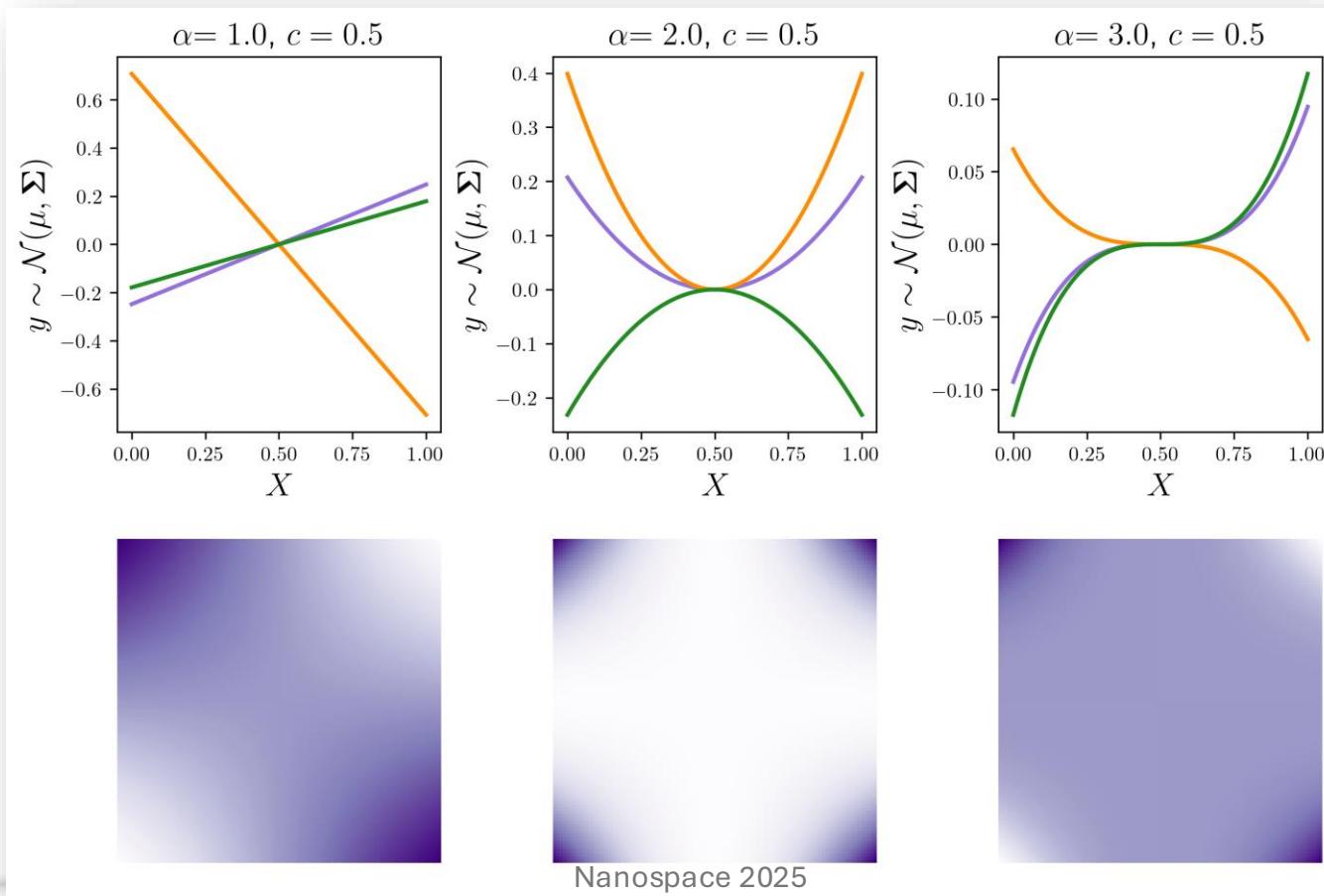


Covariance matrix → Covariance function = Kernel function

- “A kernel is a function that defines the covariance between function values at any pair of input locations.”

$$K(x_1, x_2) = ((x_1 - c) * (x_2 - c))^{\alpha}$$

Polynomial kernel



Multivariate Gaussian distribution

- Multivariate normal distribution parameterized by the mean μ and the covariance Σ .

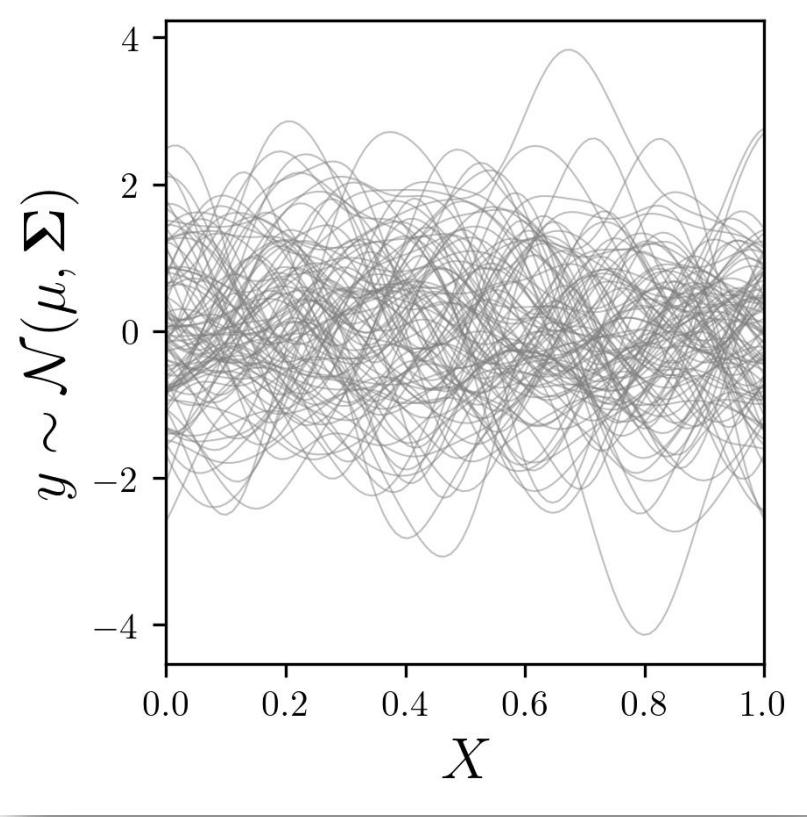
Stochastic Process:

- "A collection of random variables indexed by time, space, or some other domain."

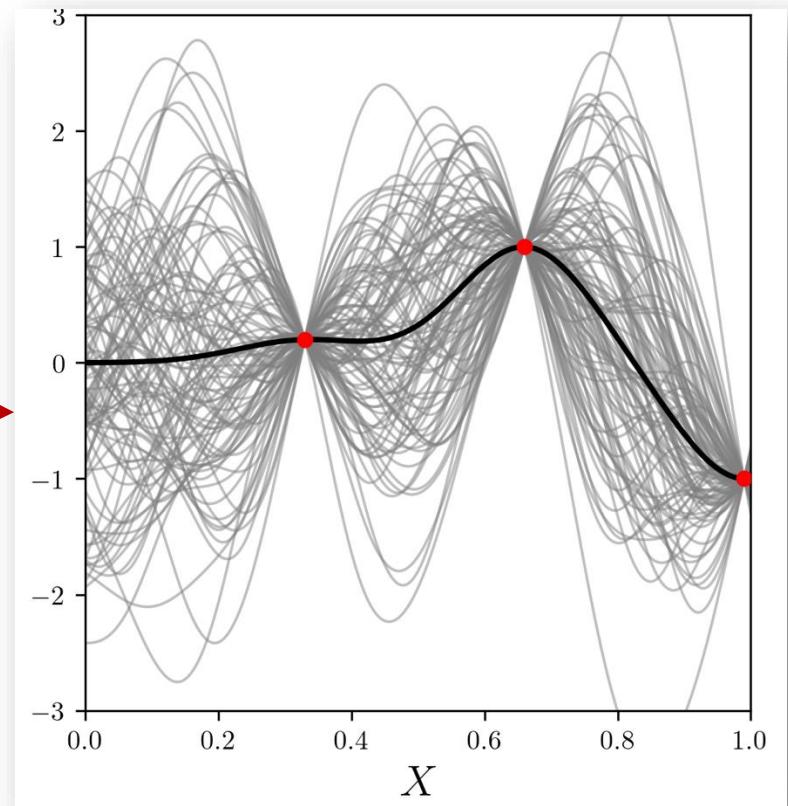
Gaussian Process

- "*A probability distribution over functions*, where any finite collection of function values is jointly *Gaussian distributed*."

Gaussian Processes & Regressors



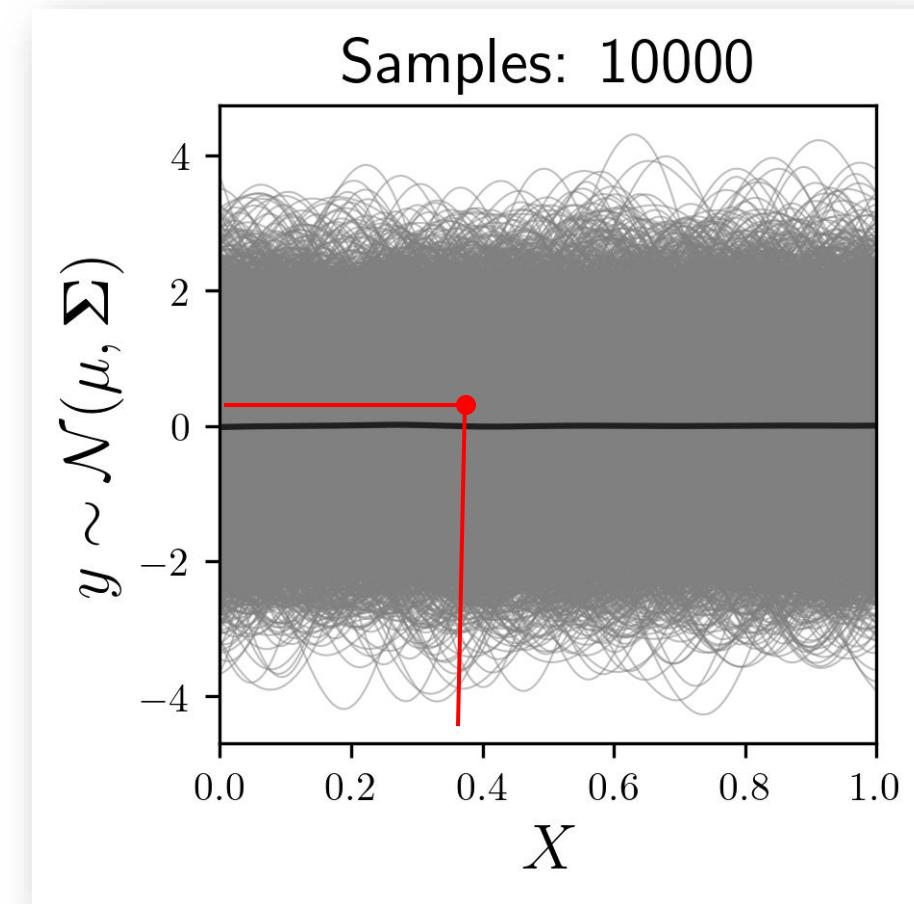
How to use for
regression?



Naïve Gaussian Process Regression:

- Given targets $f(x_{\text{target}}) = y_{\text{target}}$
- Sample a large number of functions from the GP
- Keep only those functions that are sufficiently close to y_{target} at x_{target}

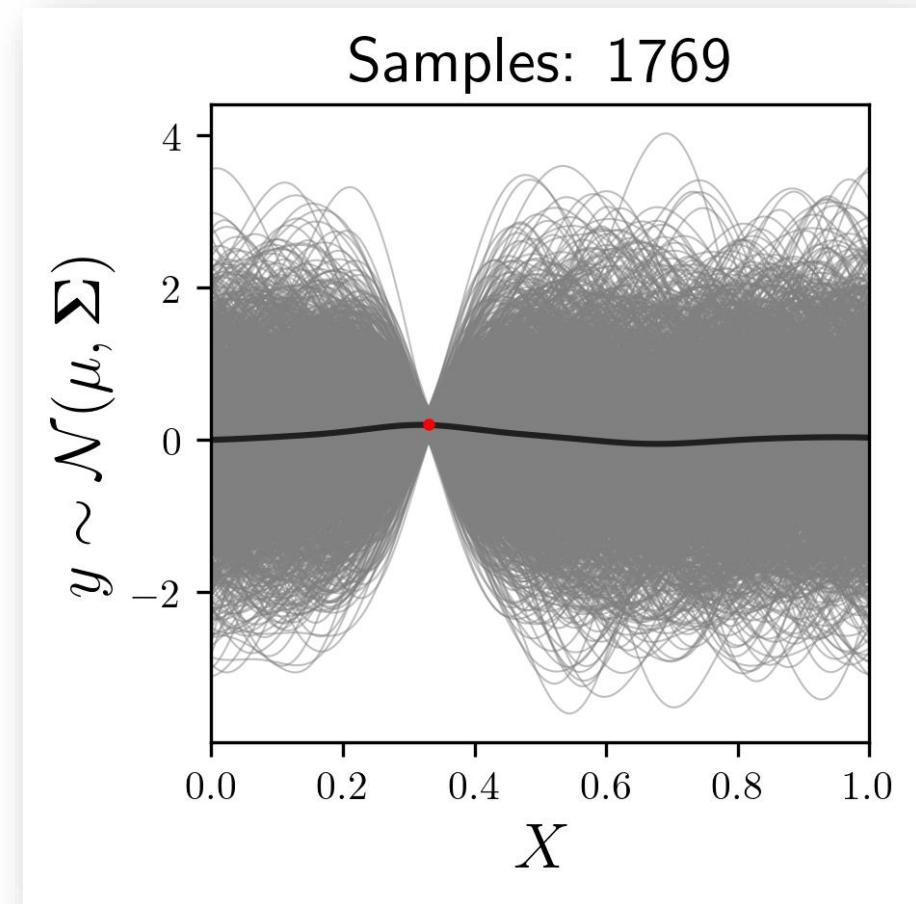
Say that our fit is the mean of the qualified samples



Naïve Gaussian Process Regression:

- Given targets $f(x_{\text{target}}) = y_{\text{target}}$
- Sample a large number of functions from the GP
- Keep only those functions that are sufficiently close to y_{target} at x_{target}

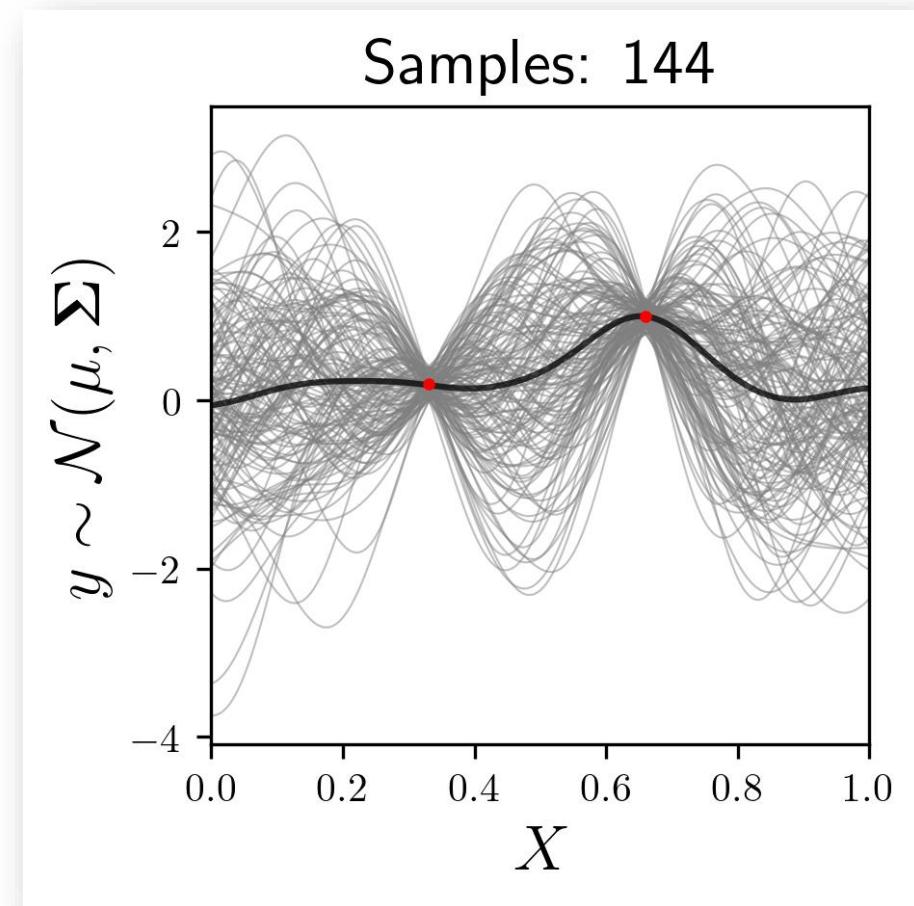
Say that our fit is the mean of the qualified samples



Naïve Gaussian Process Regression:

- Given targets $f(x_{\text{target}}) = y_{\text{target}}$
- Sample a large number of functions from the GP
- Keep only those functions that are sufficiently close to y_{target} at x_{target}

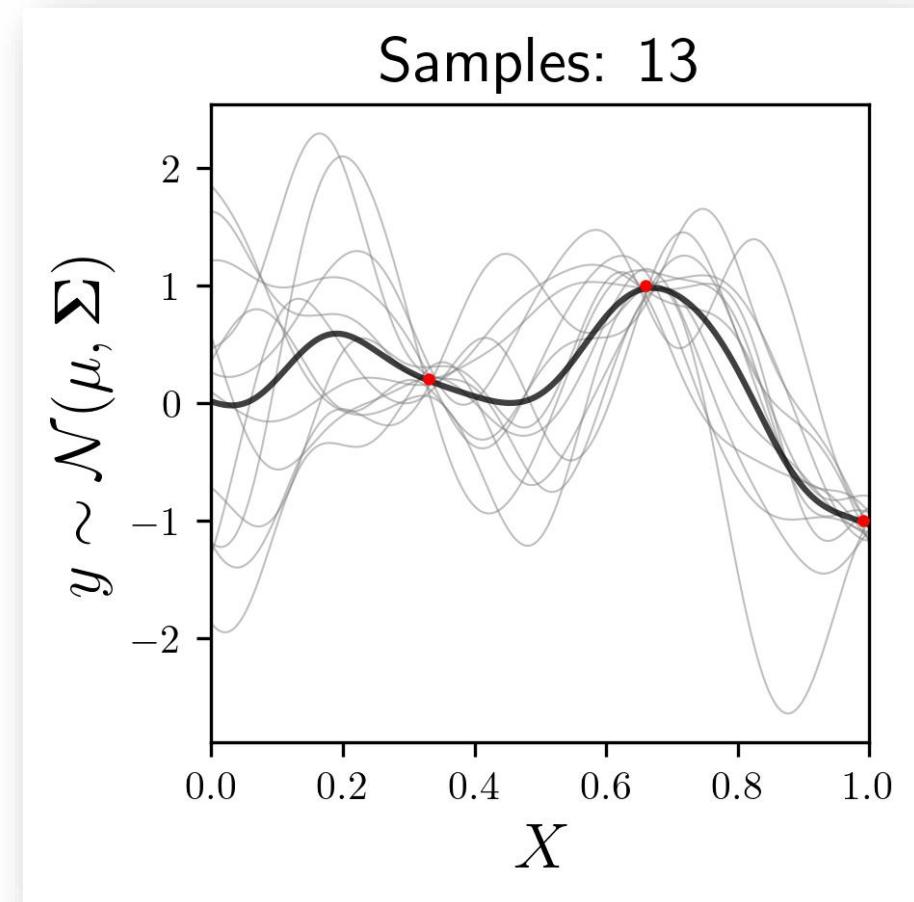
Say that our fit is the mean of the qualified samples



Naïve Gaussian Process Regression:

- Given targets $f(x_{\text{target}}) = y_{\text{target}}$
- Sample a large number of functions from the GP
- Keep only those functions that are sufficiently close to y_{target} at x_{target}

Say that our fit is the mean of the qualified samples



Naïve Gaussian Process Regression:

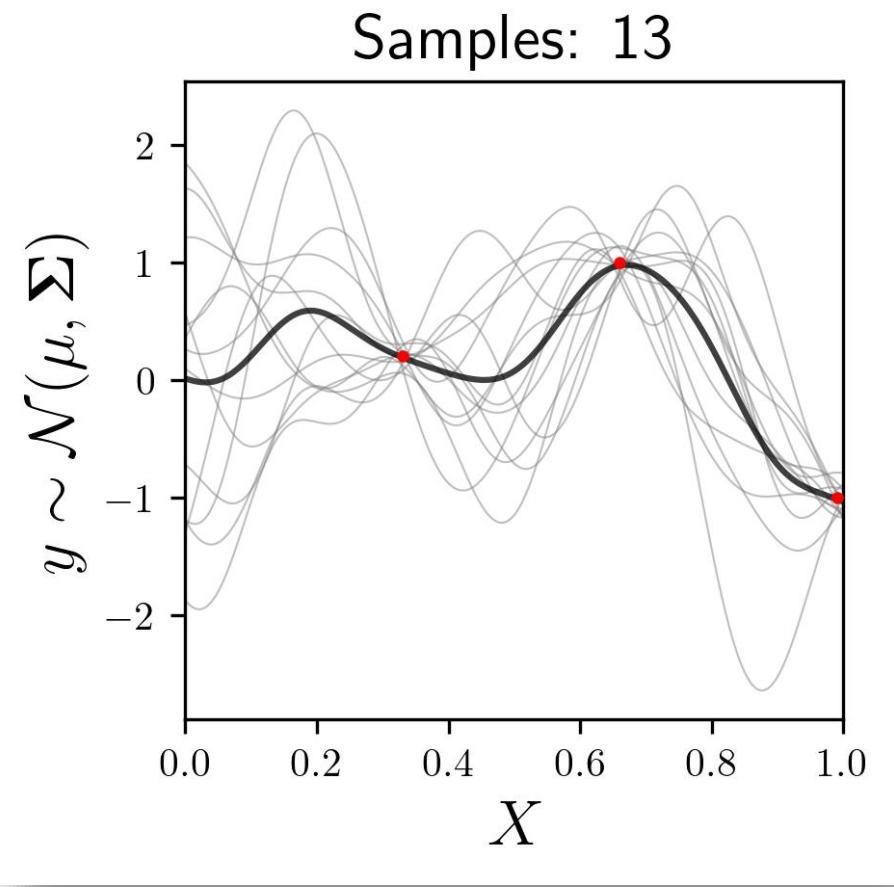
- Given targets $f(x_{\text{target}}) = y_{\text{target}}$
- Sample a large number of functions from the GP
- Keep only those functions that are sufficiently close to y_{target} at x_{target}

Say that our fit is the mean of the qualified samples

Problem: The number of samples that fit falls by about a factor of 10 pr. Added data point.

$$10000 \rightarrow 1769 \rightarrow 144 \rightarrow 13$$

Among other problems.

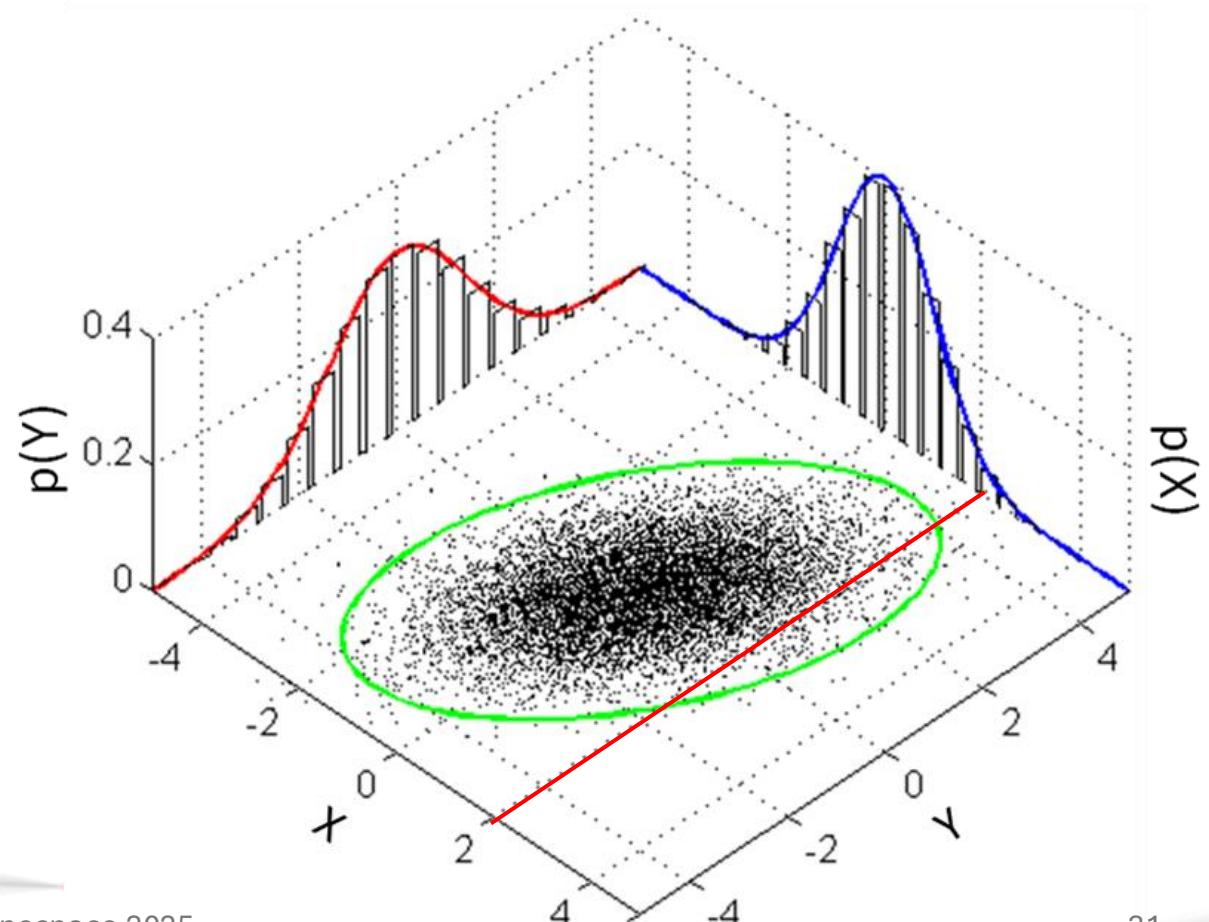


Can something smarter be done?

Bayes Theorem

- Rule to compute a conditional distribution from a joint distribution.

$$p(y|X) = \frac{p(X|y)p(y)}{p(X)}$$



Bayes Theorem

- Rule to compute a conditional distribution from a joint distribution.

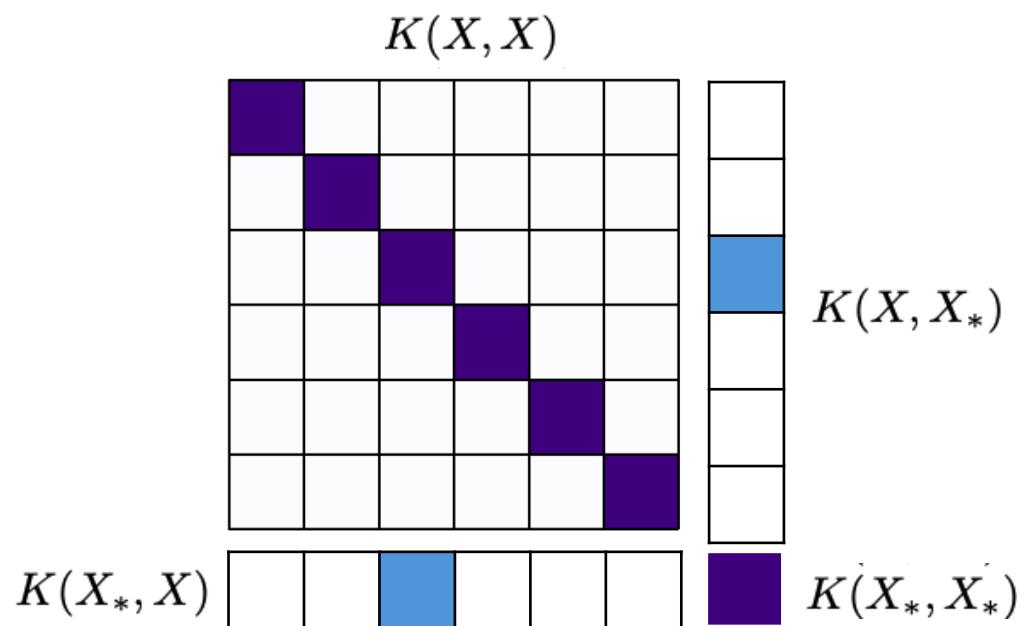
Joint prior distribution

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right).$$

Due to the ~~magical~~ properties of Gaussians the mean and covariance of the conditional distribution are

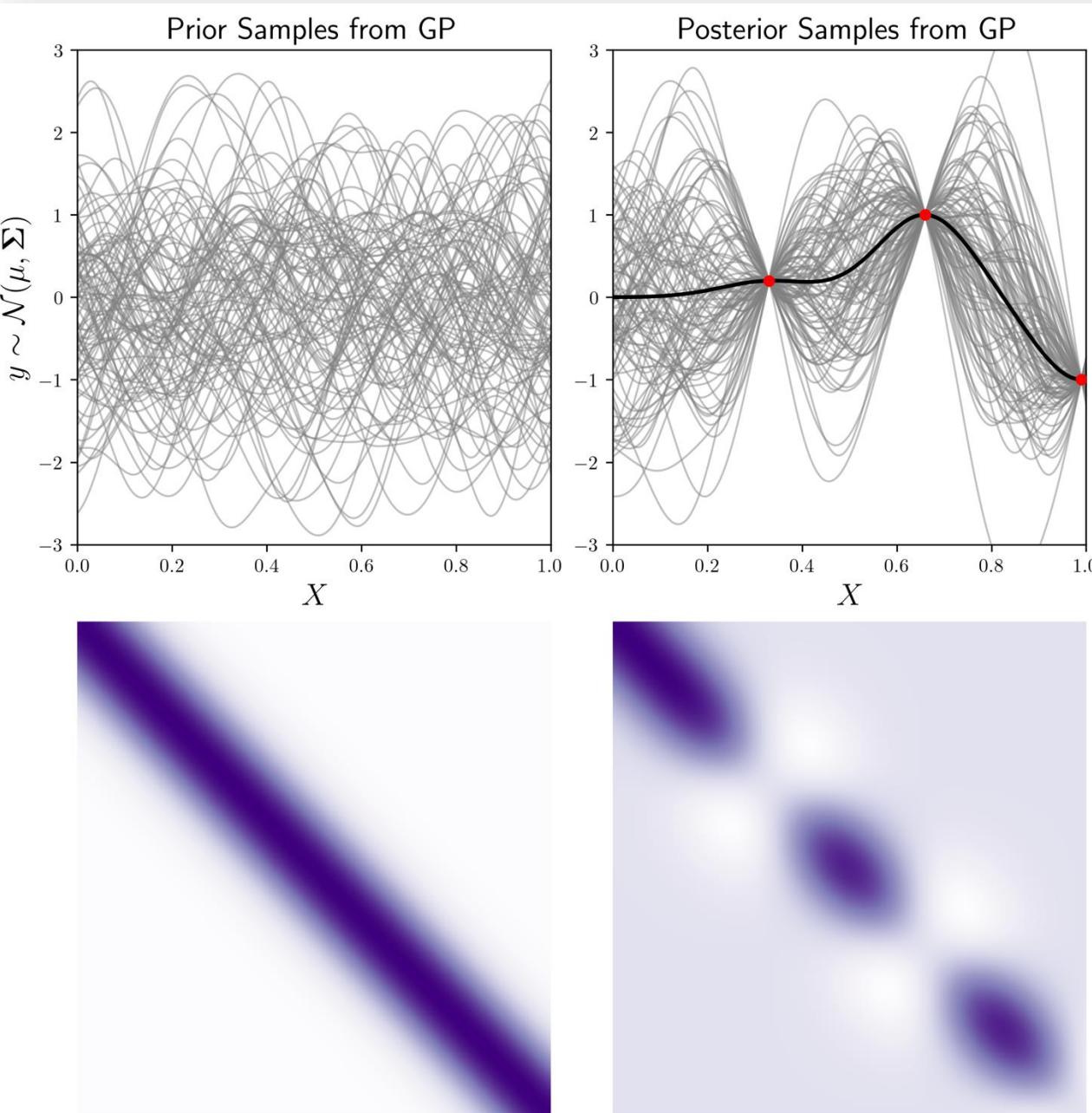
$$\boldsymbol{\mu}^* = K(X^*, X)K(X, X)^{-1}\mathbf{f}$$

$$\boldsymbol{\Sigma}^* = K(X^*, X^*) - K(X^*, X)K(X, X)^{-1}K(X, X^*)$$



In GP terminology these are typically called the **posterior** mean and covariance.

$$\mu = 0$$



$$\boldsymbol{\mu}^* = K(X^*, X)K(X, X)^{-1}\mathbf{f}$$

$$\begin{aligned}\Sigma^* \\ &= K(X^*, X^*) \\ &- K(X^*, X)K(X, X)^{-1}K(X, X^*)\end{aligned}$$

$$\boldsymbol{\mu}^* = K(X^*, X)K(X, X)^{-1}\boldsymbol{f}$$

$$K(X^*, X) \quad \quad \quad K(X, X)^{-1} \quad \quad \boldsymbol{f}$$

$$= K(X^*, X)\boldsymbol{\alpha}$$

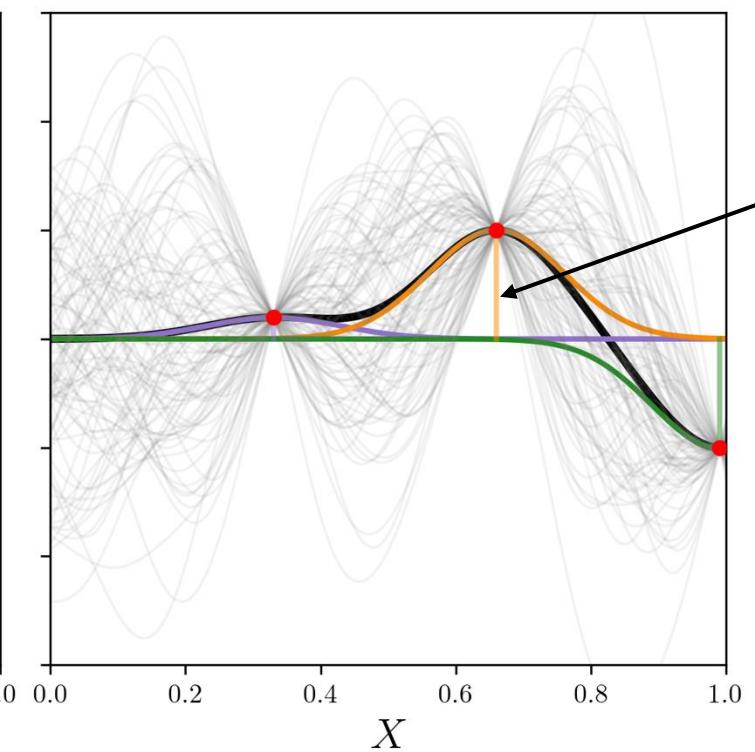
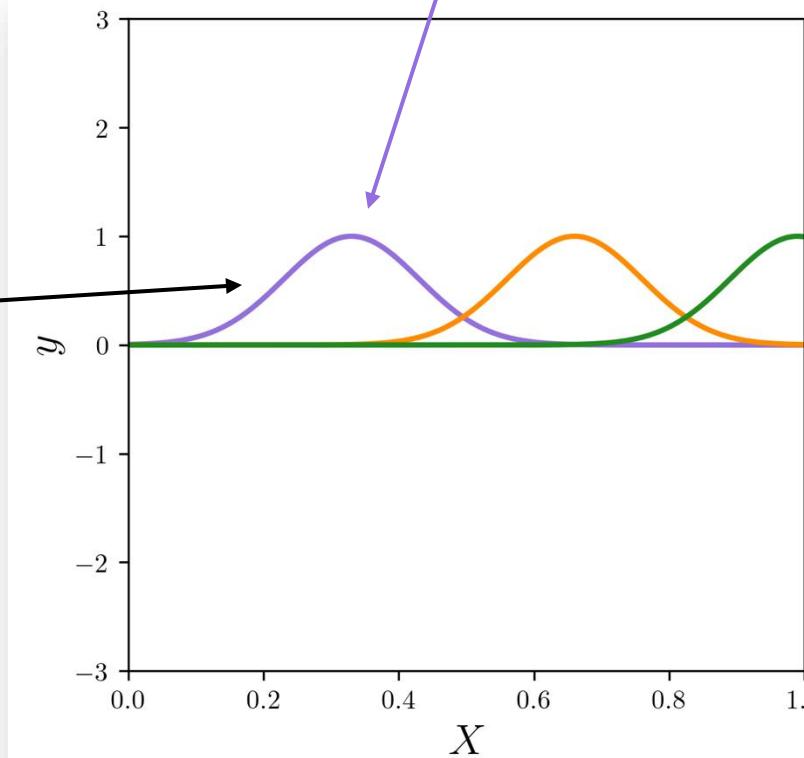
$$K(X^*, X) \quad \quad \quad \boldsymbol{\alpha}$$

$$= \sum_i \alpha_i K(X^*, X_i)$$

$$\quad \quad \quad + \quad \quad \quad + \quad \quad \quad$$

$$\mu^* = \sum_i \alpha_i K(X^*, X_i) = \text{[pink square]} \text{ [orange square]} + \text{[purple square]} \text{ [brown square]} + \text{[pink square]} \text{ [brown square]}$$

Kernels placed at observations as basis functions



Weights α_i come from the conditioning

Bayes Theorem

- Rule to compute a conditional distribution from a joint distribution.

Joint prior distribution with noisy observations

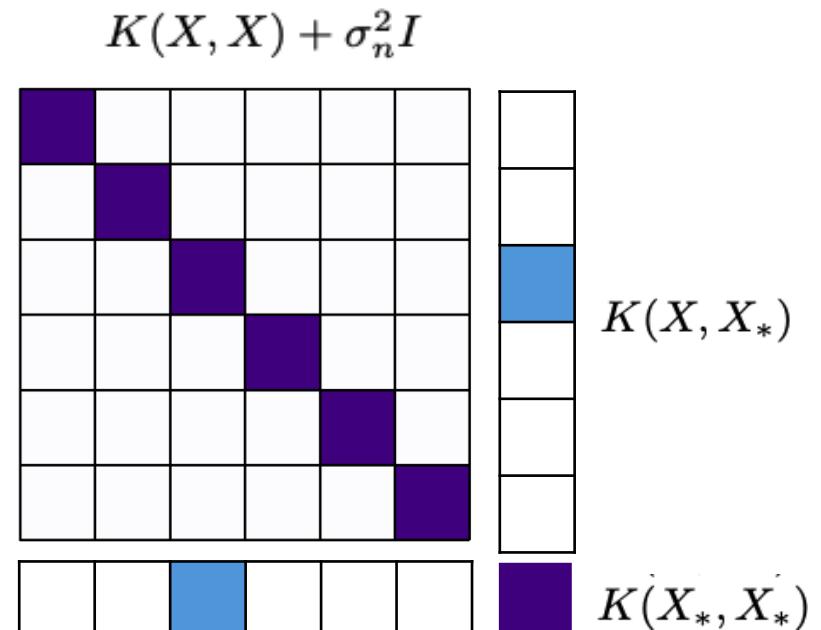
$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right).$$

Due to the magical properties of Gaussians the mean and covariance of the conditional distribution are

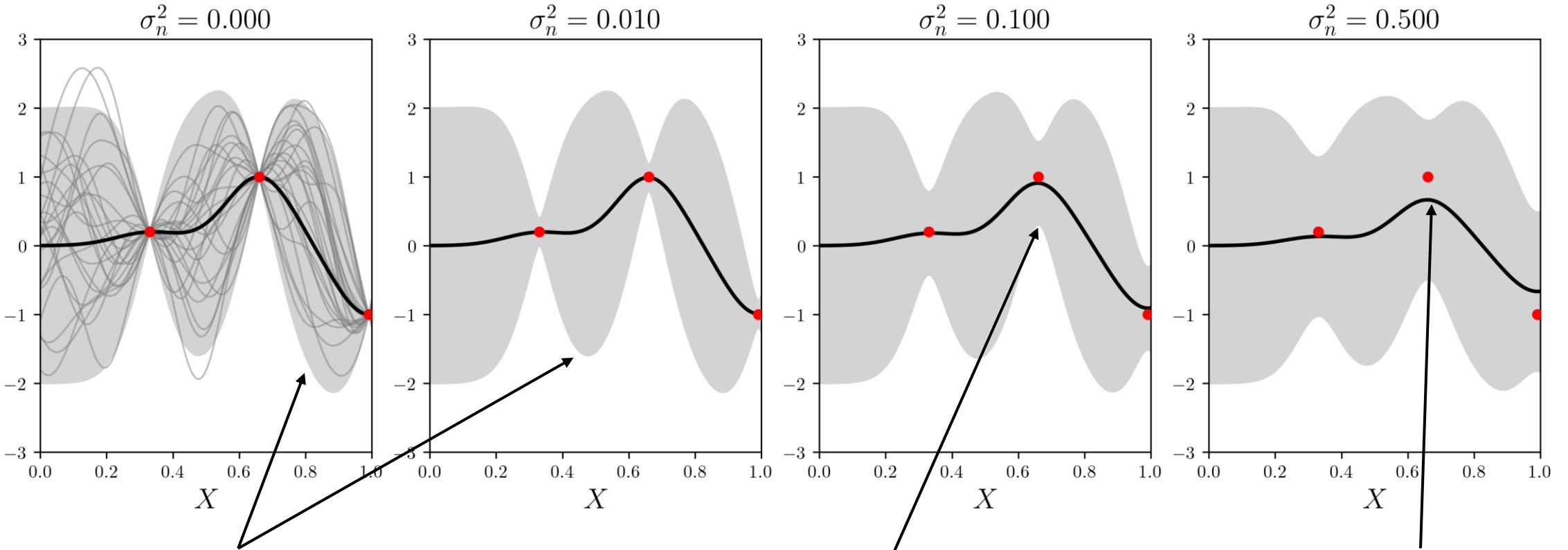
$$\boldsymbol{\mu}^* = K(X^*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{f}$$

$$\boldsymbol{\Sigma}^* = K(X^*, X^*) - K(X^*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X^*)$$

In GP terminology these are typically called the **posterior** mean and covariance.



$$\boldsymbol{\mu}^* = K(X^*, X)[K(X, X)^{-1} + \sigma_n^2 I]\mathbf{f}$$



Rather than plotting samples
we can just plot the **2 std**
confidence interval resulting
from the covariance.

More noisy
observations →
Less confident

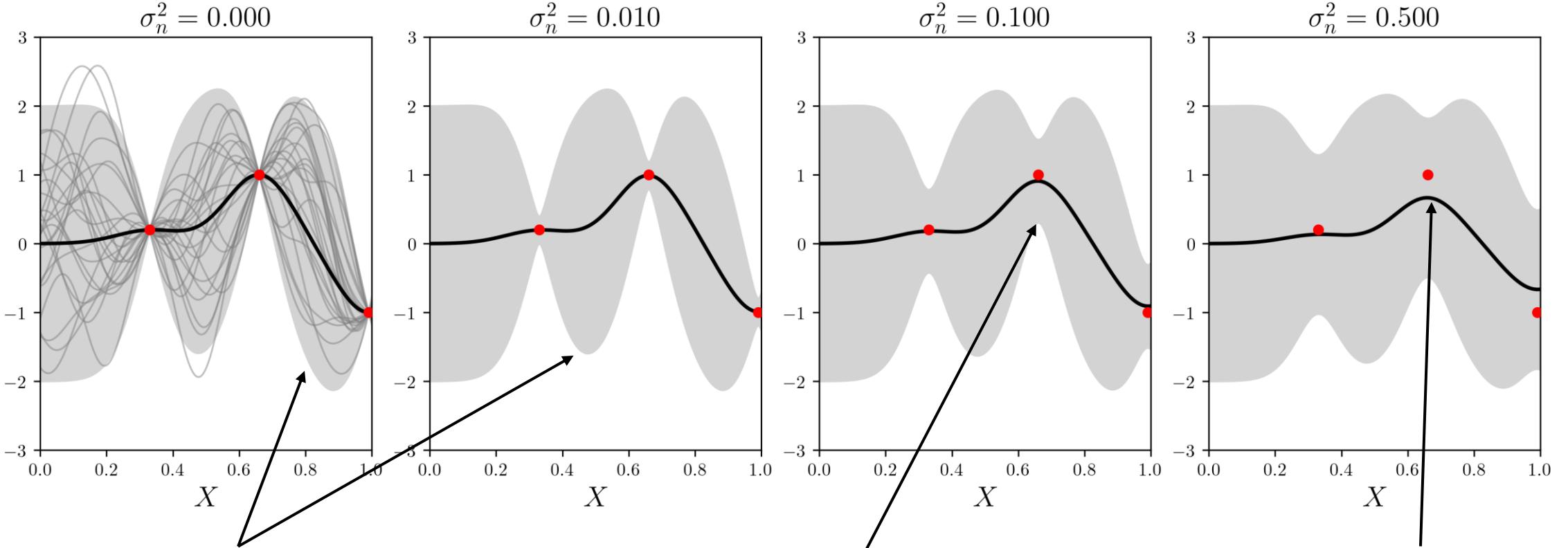
Posterior mean less
obliged to hit the data →
Reverts to prior mean = 0

Kernel ridge regression

$$\begin{aligned}\mu^* &= K(X^*, X)[K(X, X)^{-1} + \sigma_n^2 I]f \\ &= \sum_i \alpha_i K(X^*, X_i)\end{aligned}$$



$$L(\alpha) = \left\| \sum_i \alpha_i K(X^*, X_i) - y_i \right\|^2 + \sigma_n^2 \|\alpha\|^2$$



Rather than plotting samples
we can just plot the **2 std
confidence interval** resulting
from the covariance.

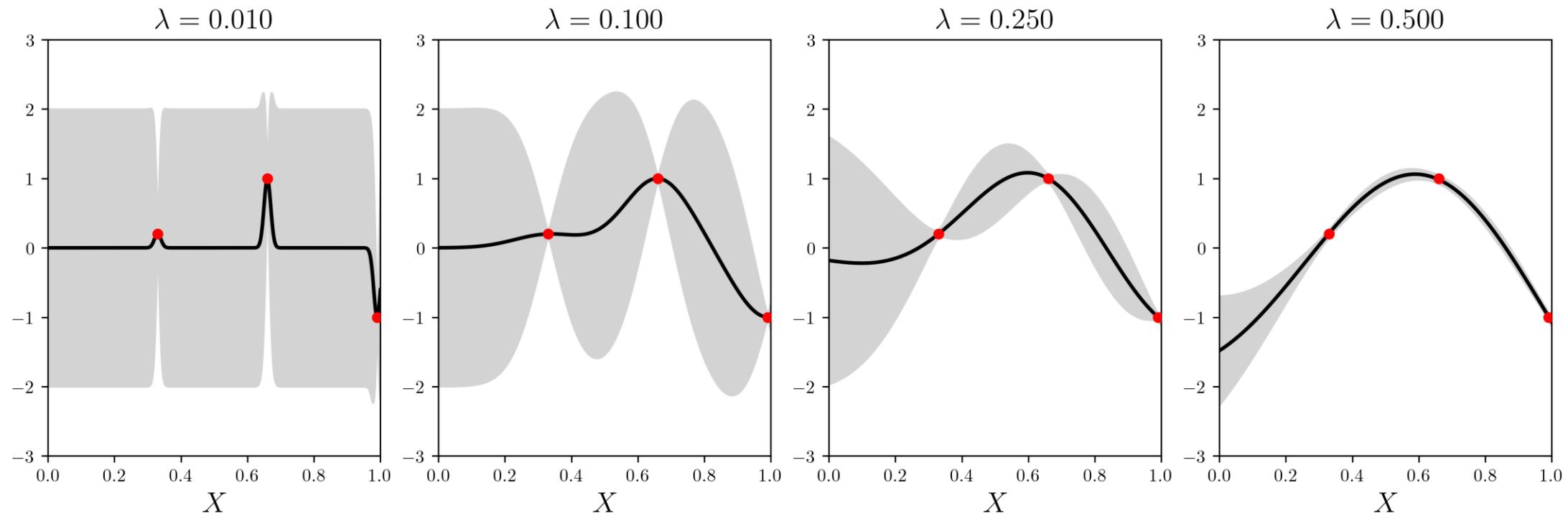
More noisy
observations →
Less confident

Posterior mean less
obliged to hit the data →
Reverts to prior mean = 0

Covariance matrix → Covariance function = Kernel function

- “A kernel is a function that defines the covariance between function values at any pair of input locations.”

$$K(x_1, x_2) = \exp\left(-\frac{\|x_2 - x_1\|^2}{2\lambda^2}\right)$$

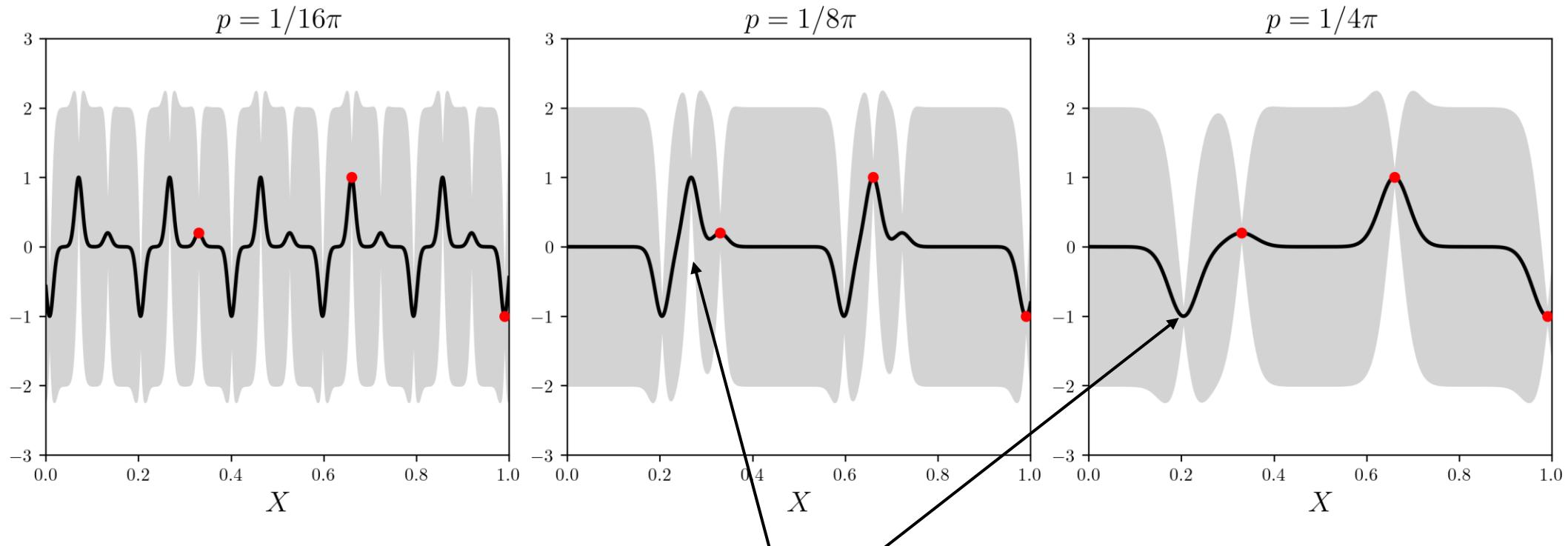


Longer length scale λ gives a mean that changes slowly and the “relative distance” to observation is smaller so less variance → Tighter confidence intervals.

Covariance matrix → Covariance function = Kernel function

- “A kernel is a function that defines the covariance between function values at any pair of input locations.”

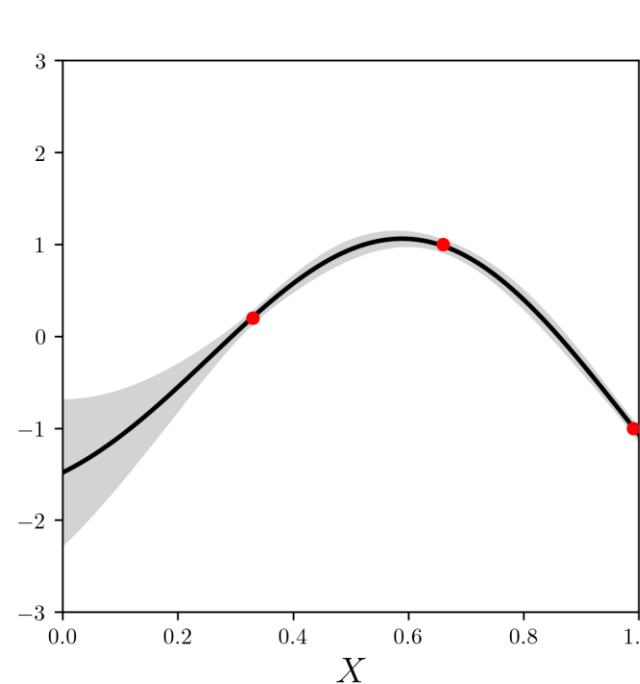
$$K(x_1, x_2) = \exp\left(-\frac{2}{\sigma^2} \sin^2\left(\frac{\pi \|x_1 - x_2\|}{p}\right)\right) \quad \text{Periodic kernel}$$



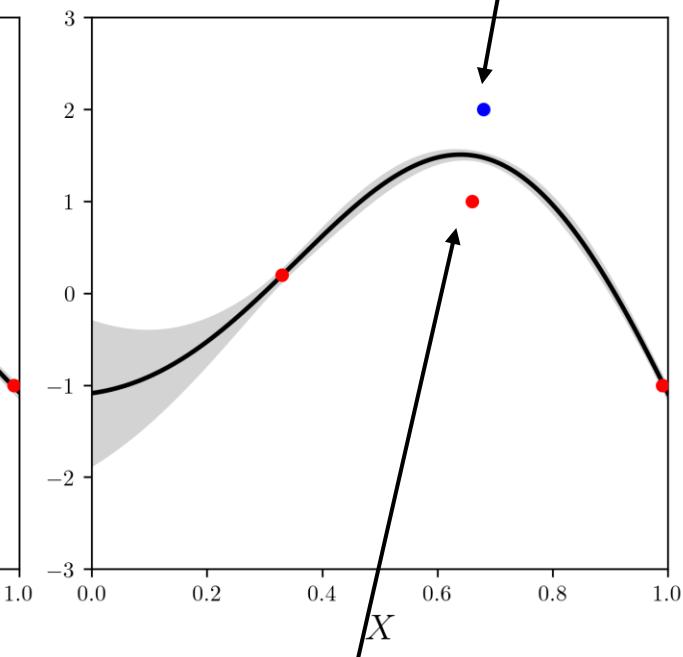
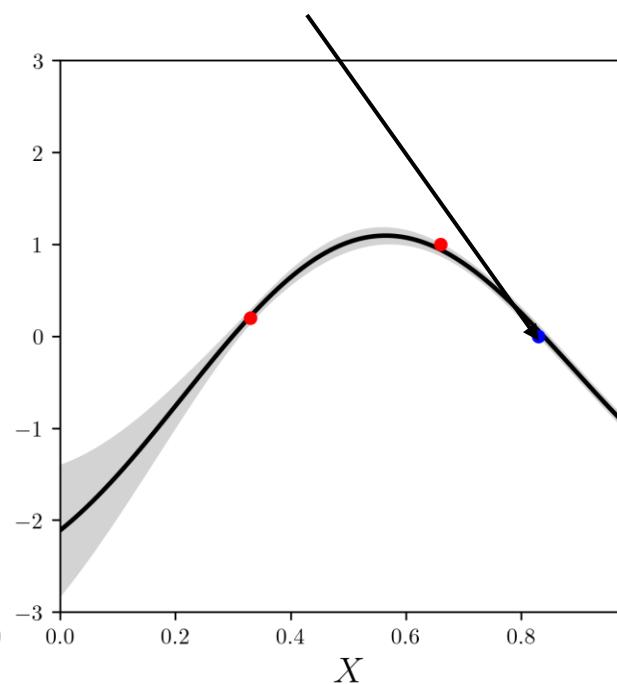
The periodic kernel introduces low variance regions where we don't have direct observations

$$K(x_1, x_2) = \exp\left(-\frac{\|x_2 - x_1\|^2}{2\lambda^2}\right)$$

Additional datapoint fits well



Additional datapoint fits poorly



Distribution does not contain functions that can fit all the points

The kernel function has poor parameters!

Log marginal likelihood

- “Log-probability of observing the data under the GP model, integrating over all possible latent functions”

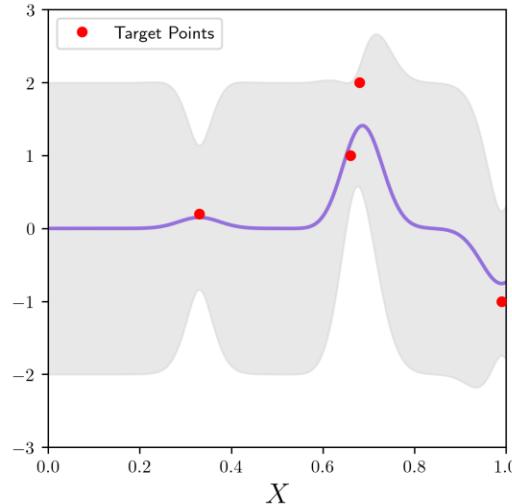
$$p(\mathbf{y}|X) = \int p(\mathbf{y}|\mathbf{f}, X)p(\mathbf{f}|X) d\mathbf{f}.$$

Again Gaussians have magical nice properties

$$\log p(\mathbf{y}|X) = -\frac{1}{2} \mathbf{y}^T (K_\theta(X, X) + \sigma_n^2 I)^{-1} \mathbf{y} \quad \text{Data fit}$$

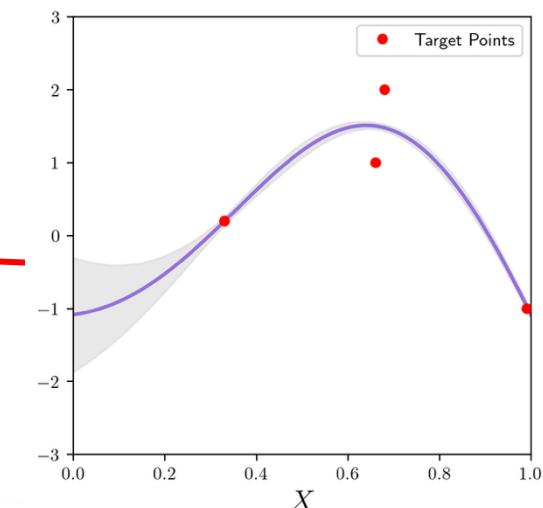
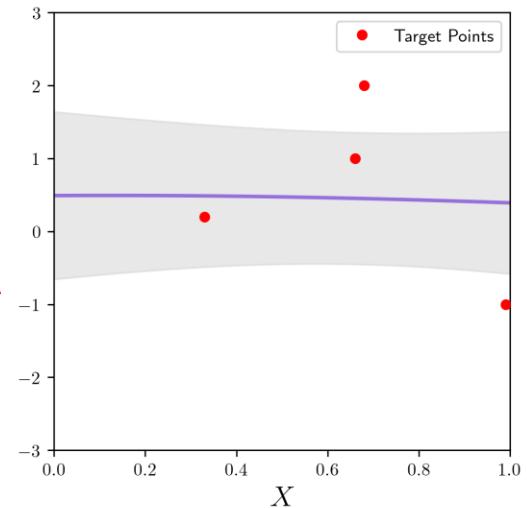
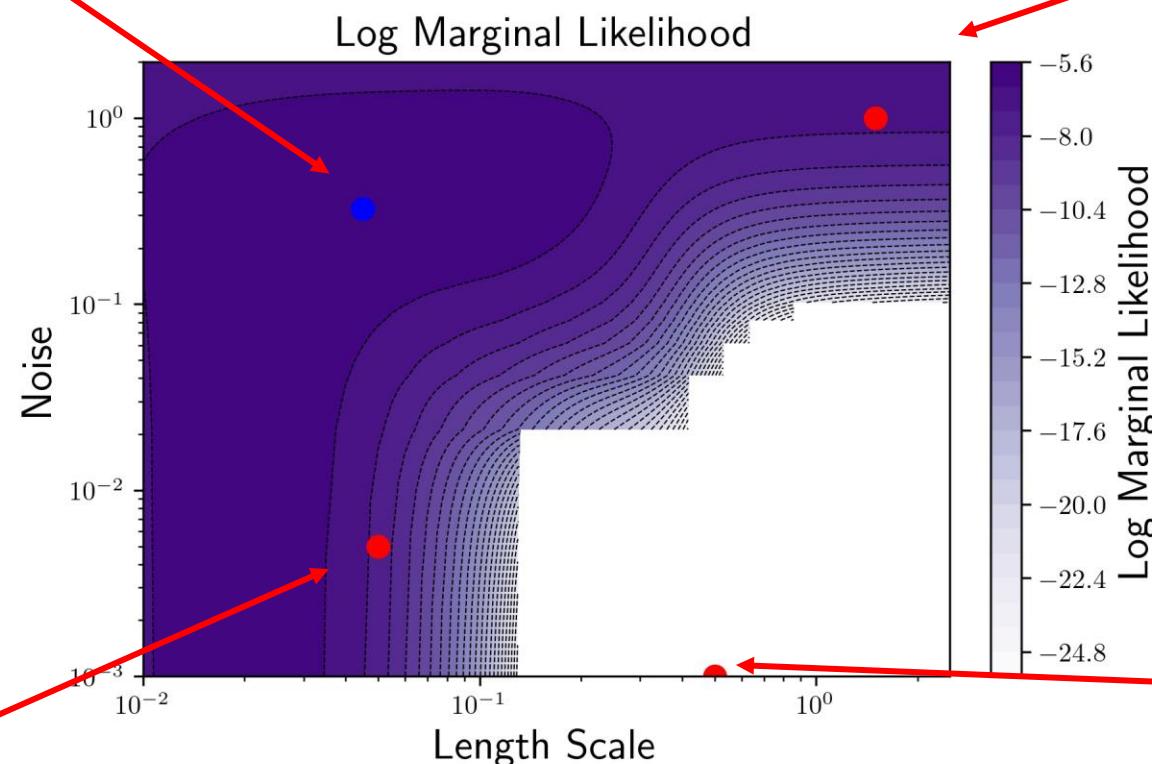
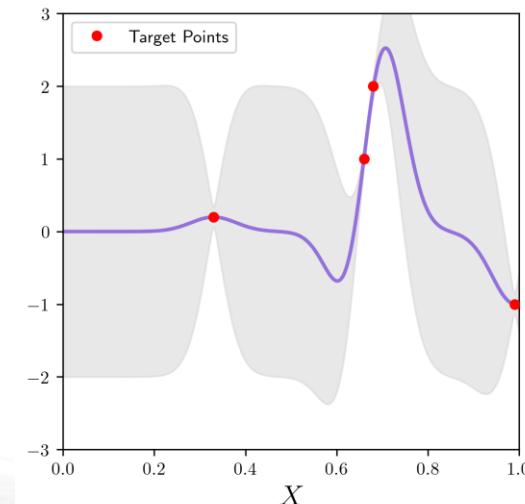
$$-\frac{1}{2} \log |K_\theta(X, X) + \sigma^2 I| \quad \text{Model complexity}$$

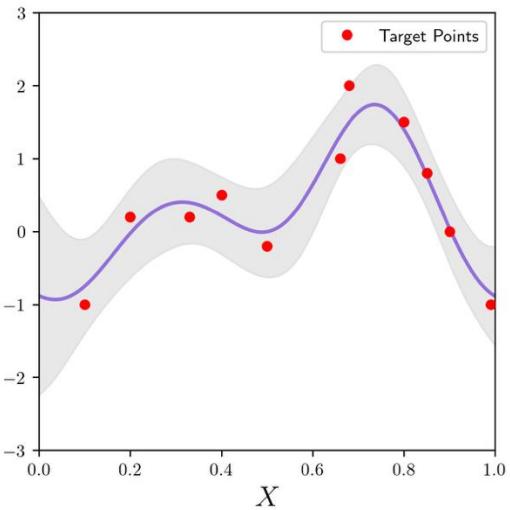
$$-\frac{n}{2} \log 2\pi \quad \text{Normalization}$$



$$\log p(\mathbf{y}|X) = -\frac{1}{2} \mathbf{y}^T (K_\theta(X, X) + \sigma_n^2 I)^{-1} \mathbf{y}$$

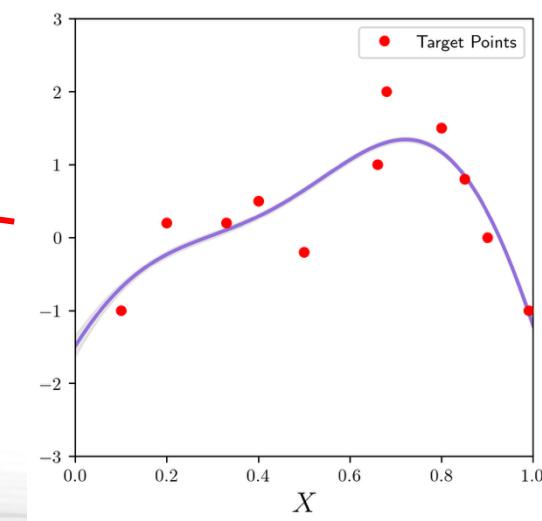
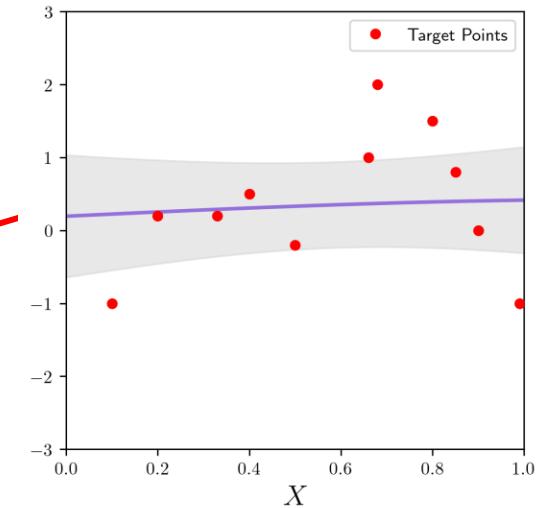
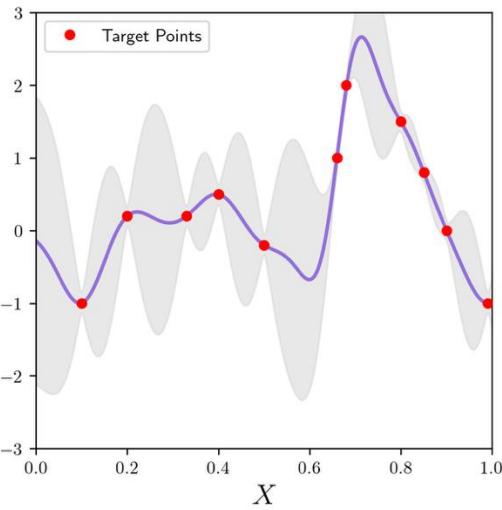
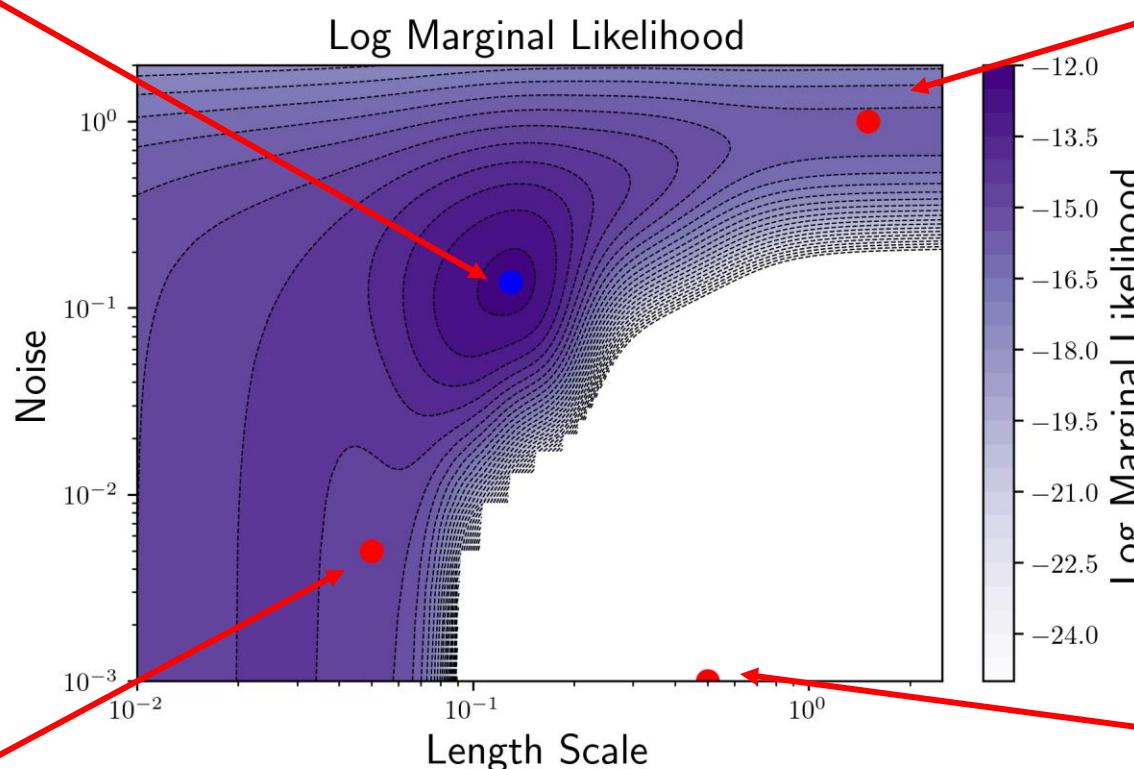
$$-\frac{1}{2} \log|K_\theta(X, X) + \sigma^2 I| - \frac{n}{2} \log 2\pi$$

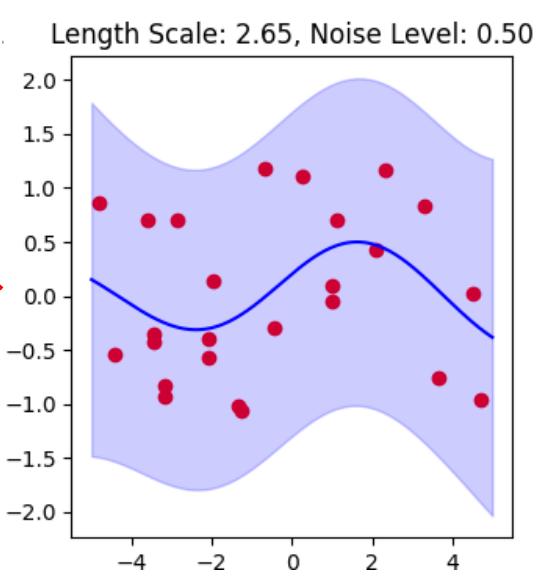
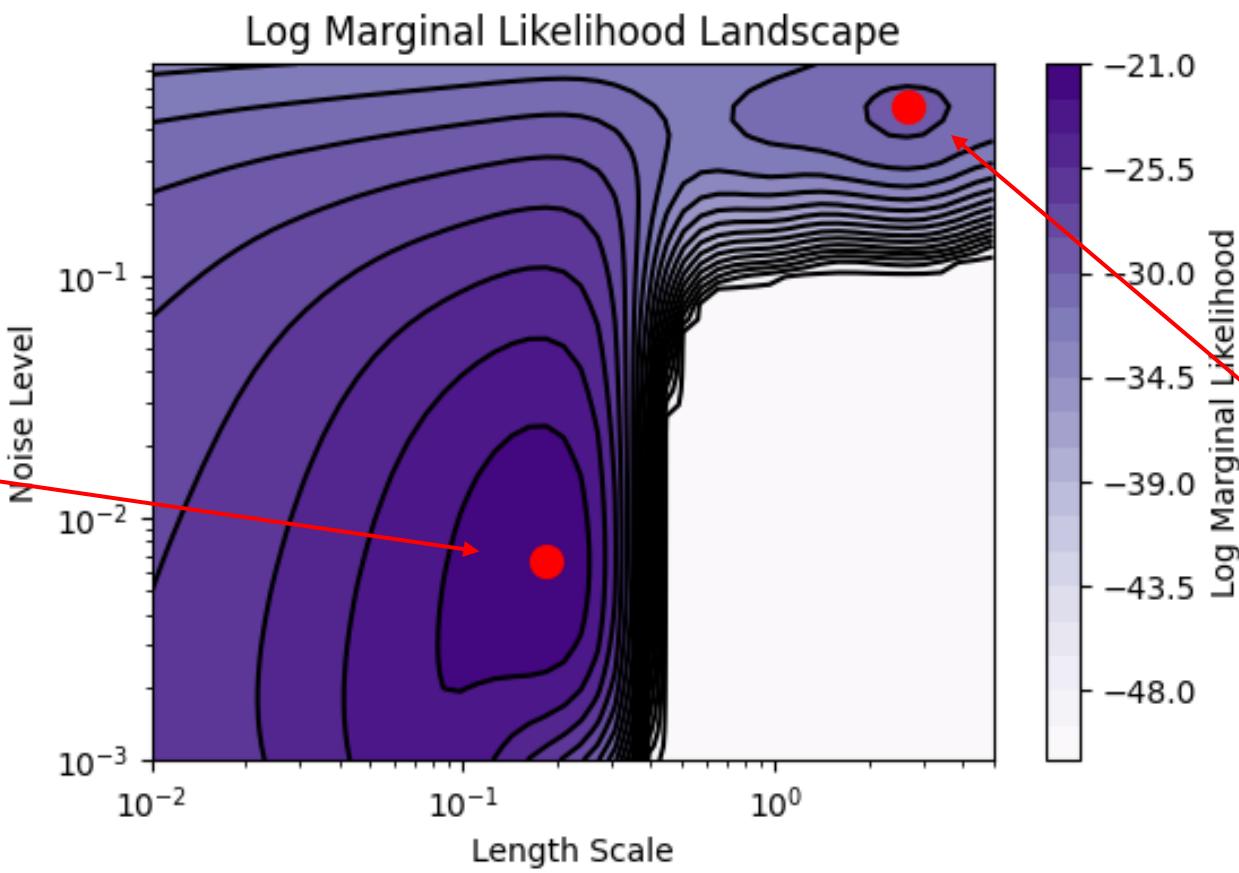
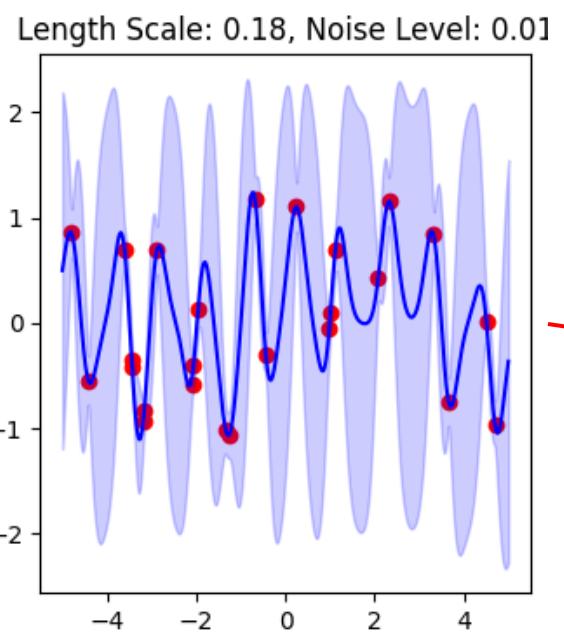




$$\log p(\mathbf{y}|X) = -\frac{1}{2} \mathbf{y}^T (K_\theta(X, X) + \sigma_n^2 I)^{-1} \mathbf{y}$$

$$-\frac{1}{2} \log|K_\theta(X, X) + \sigma^2 I| - \frac{n}{2} \log 2\pi$$

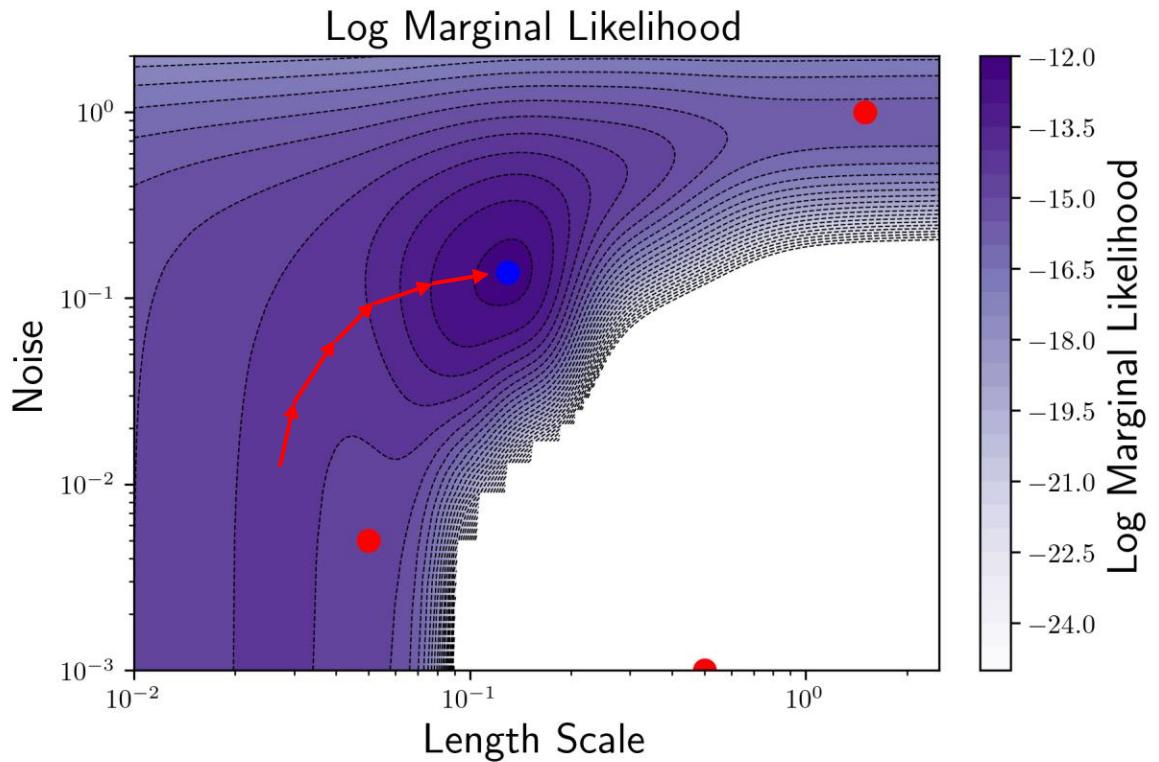




LML typically optimized using gradient-based approach

$$\nabla_{\theta} \log p_{\theta}(y|X)$$

Therefore, I recommend using a auto-differentiation framework for GPs (e.g. PyTorch or Jax)



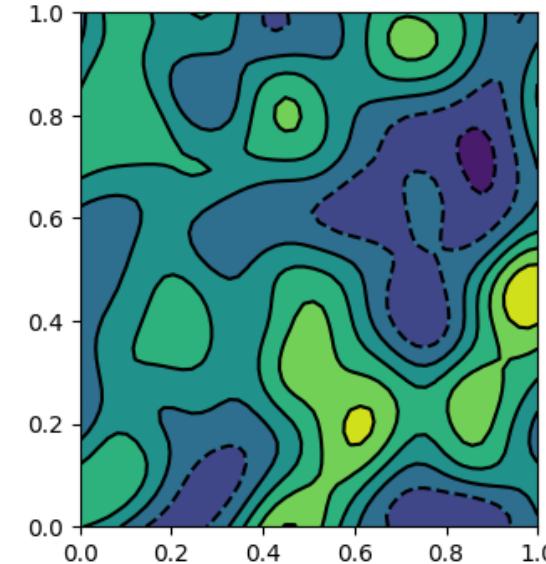
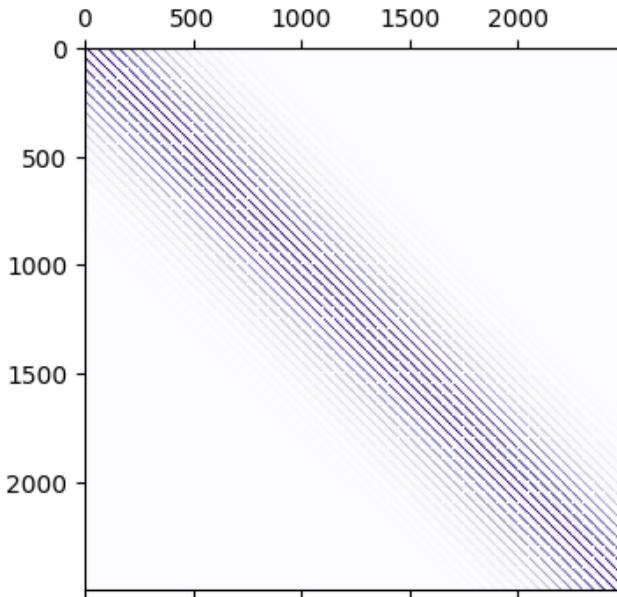
Functions of multiple variables

- All handled by the kernel – and all the kernels we've seen can do so.

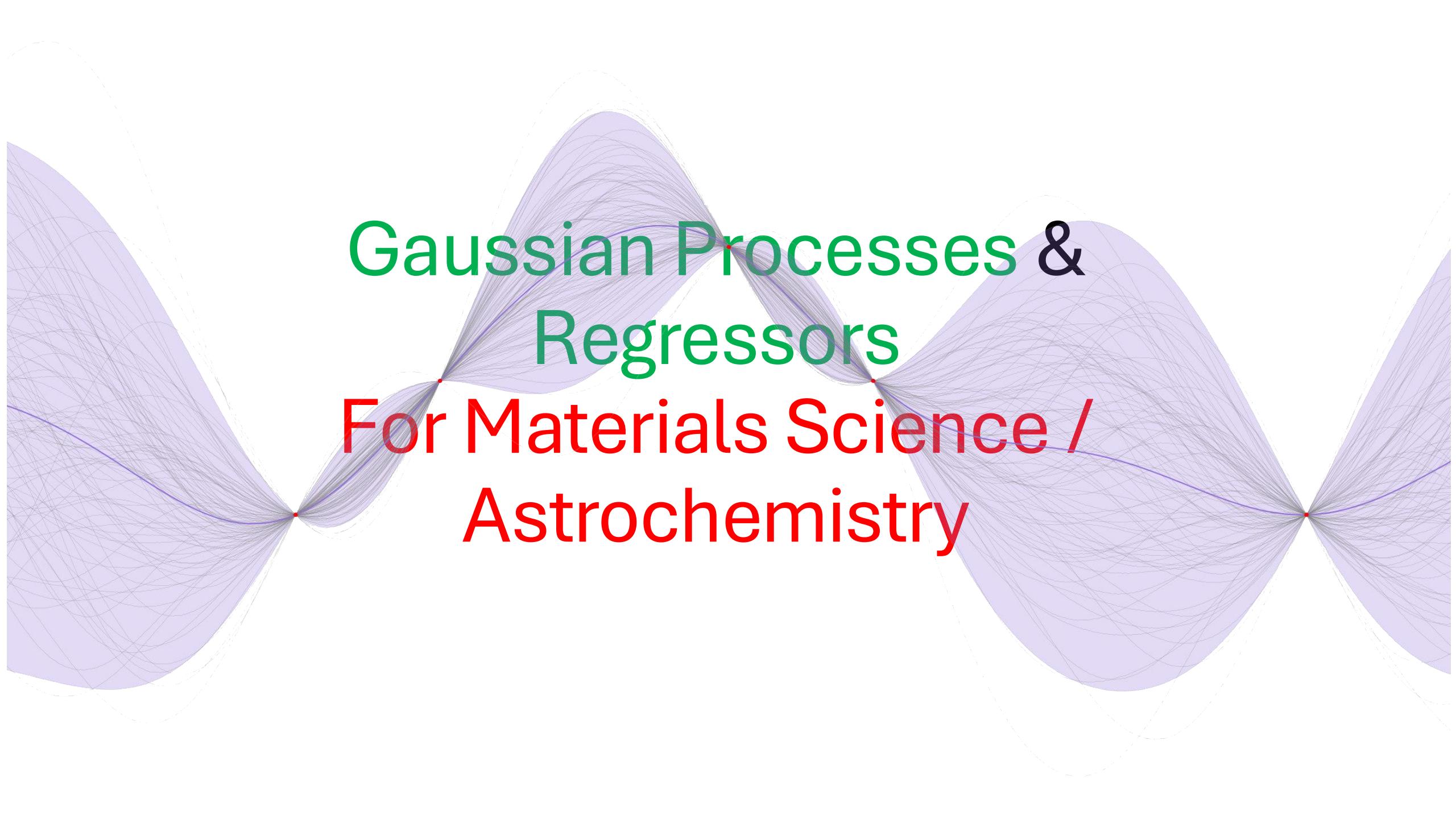
$$X = [x_1, x_2, \dots, x_n] \rightarrow X = [(x_1, y_1), (x_2, y_2), \dots (x_N, y_N)] = [\mathbf{x}_1, \mathbf{x}_2, \dots \mathbf{x}_N]$$

$$K(x_1, x_2) = \exp\left(-\frac{\|\mathbf{x}_2 - \mathbf{x}_1\|^2}{2\sigma^2}\right)$$

RBF Covariance
of a grid of points
in $[0, 1]$ for both x
and y.



Function
sampled from the
GP prior.



Gaussian Processes & Regressors For Materials Science / Astrochemistry

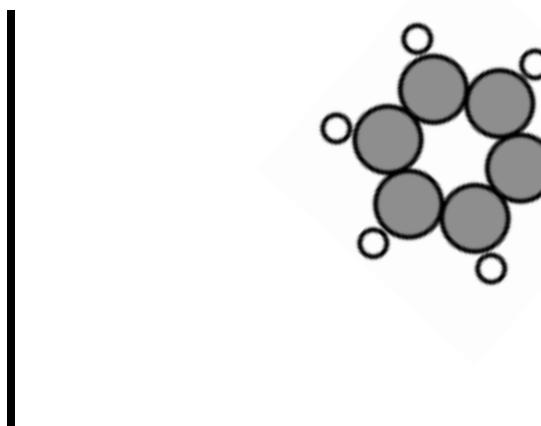
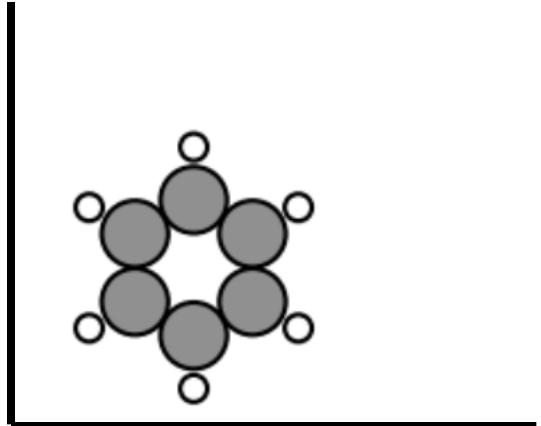
Covariance of atomic configurations.

$$K \left(\begin{array}{c} \text{atom configuration} \\ \text{graph representation} \end{array}, \begin{array}{c} \text{atom configuration} \\ \text{graph representation} \end{array} \right)$$



Structural descriptor

- A representation of an atomic structure that encodes invariances.



Invariances

- ❖ Translation
- ❖ Rotation
- ❖ Permutation

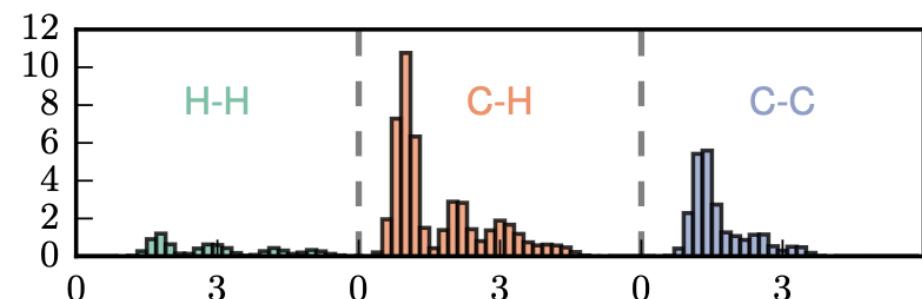
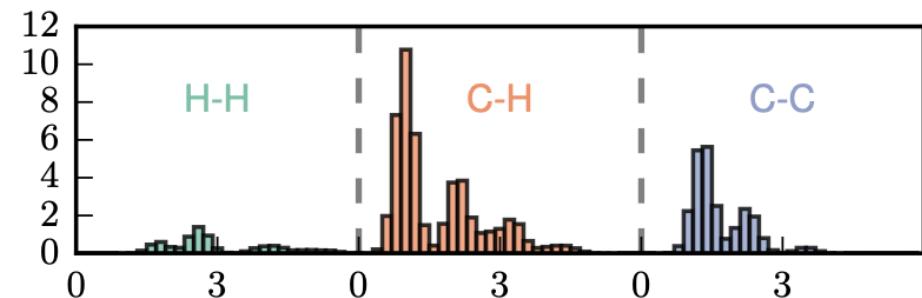
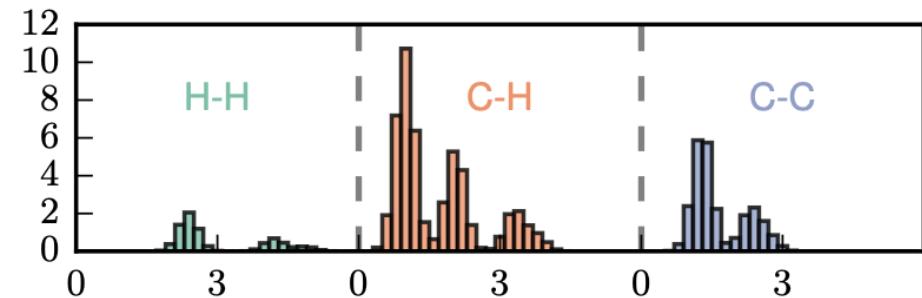
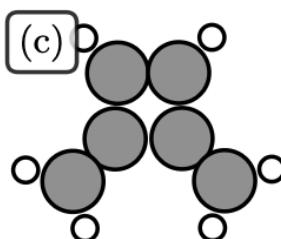
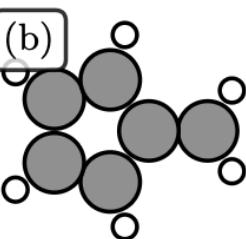
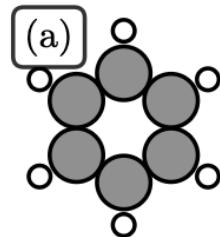
Cartesian coordinates are terrible at this.

x_1	y_1	z_1
x_2	y_2	z_2
...
x_N	y_N	z_N

Fingerprint descriptor

- A smooth histogram of bond lengths (and angles) representing an atomic configuration as a vector.

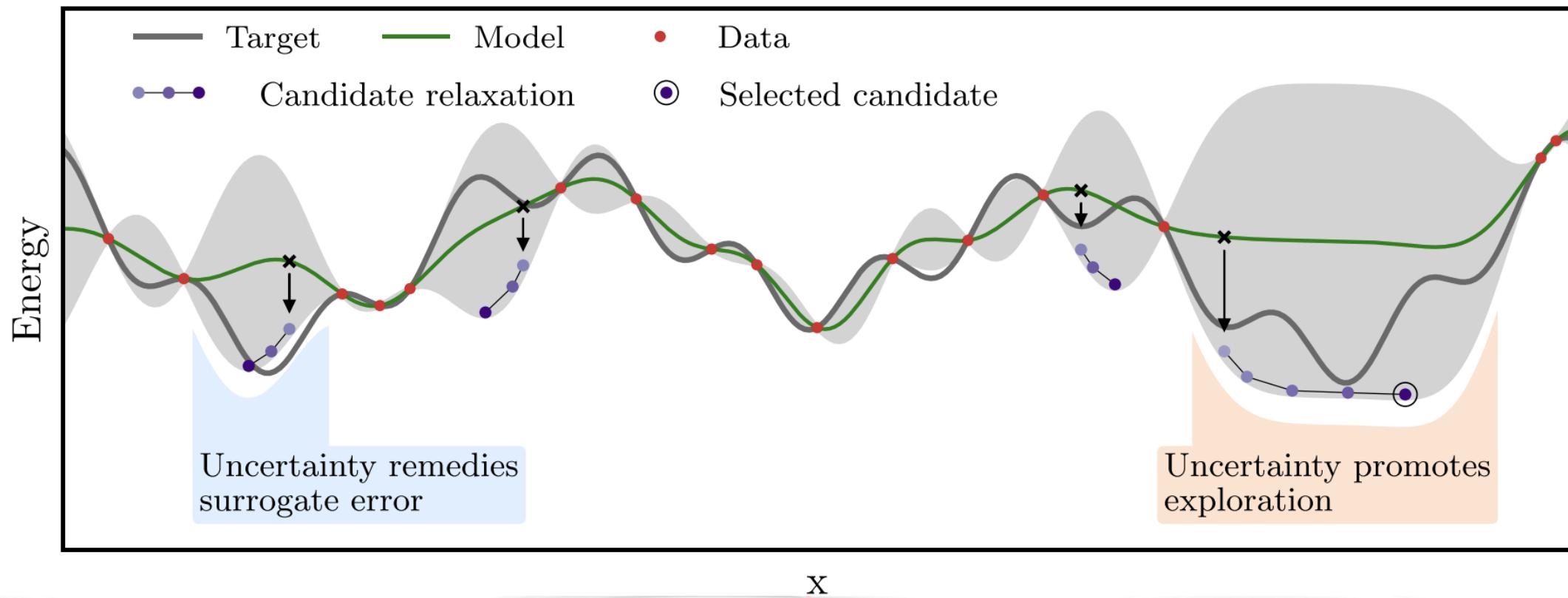
$$F_{A,B}(r) = \frac{V}{4\pi N_A N_B \Delta} \sum_i \sum_j \frac{1}{r_{ij}} \exp\left(-\frac{1}{2} \frac{(r - r_{ij}^2)}{\sigma}\right).$$

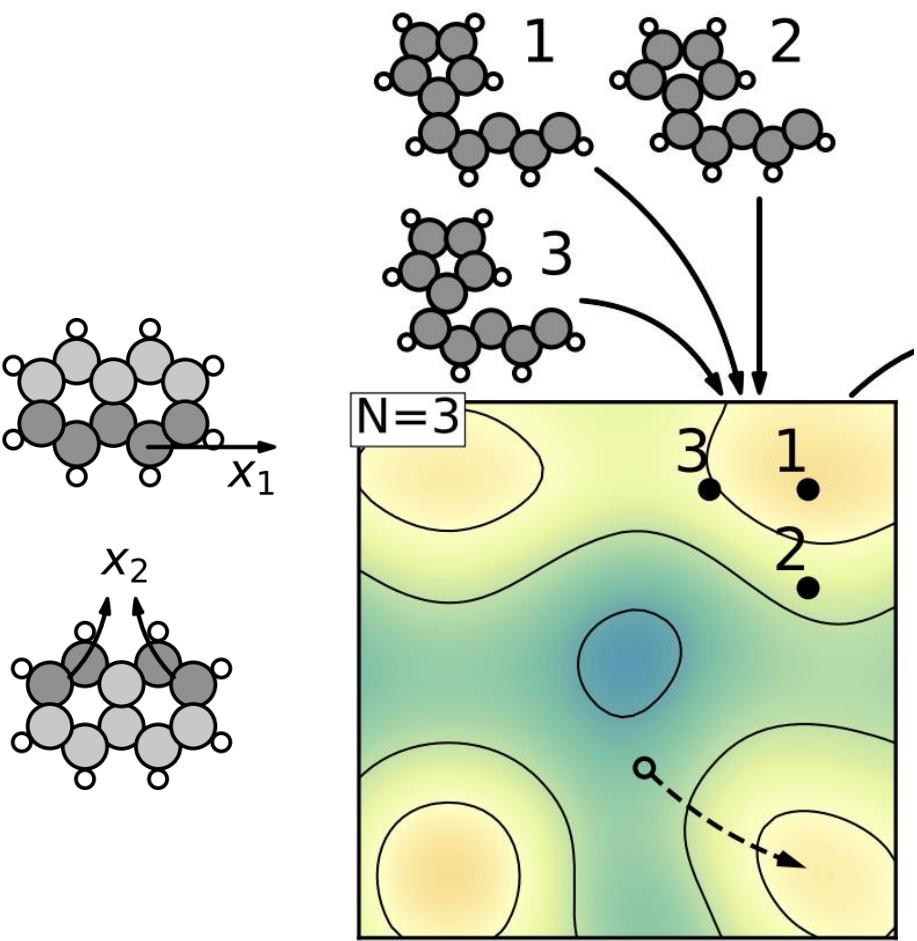


Bayesian Optimization

- Acquisition function based on the GP decides where the target function should be queried.

$$f(X) = \mu(x) - \kappa\sigma(x)$$



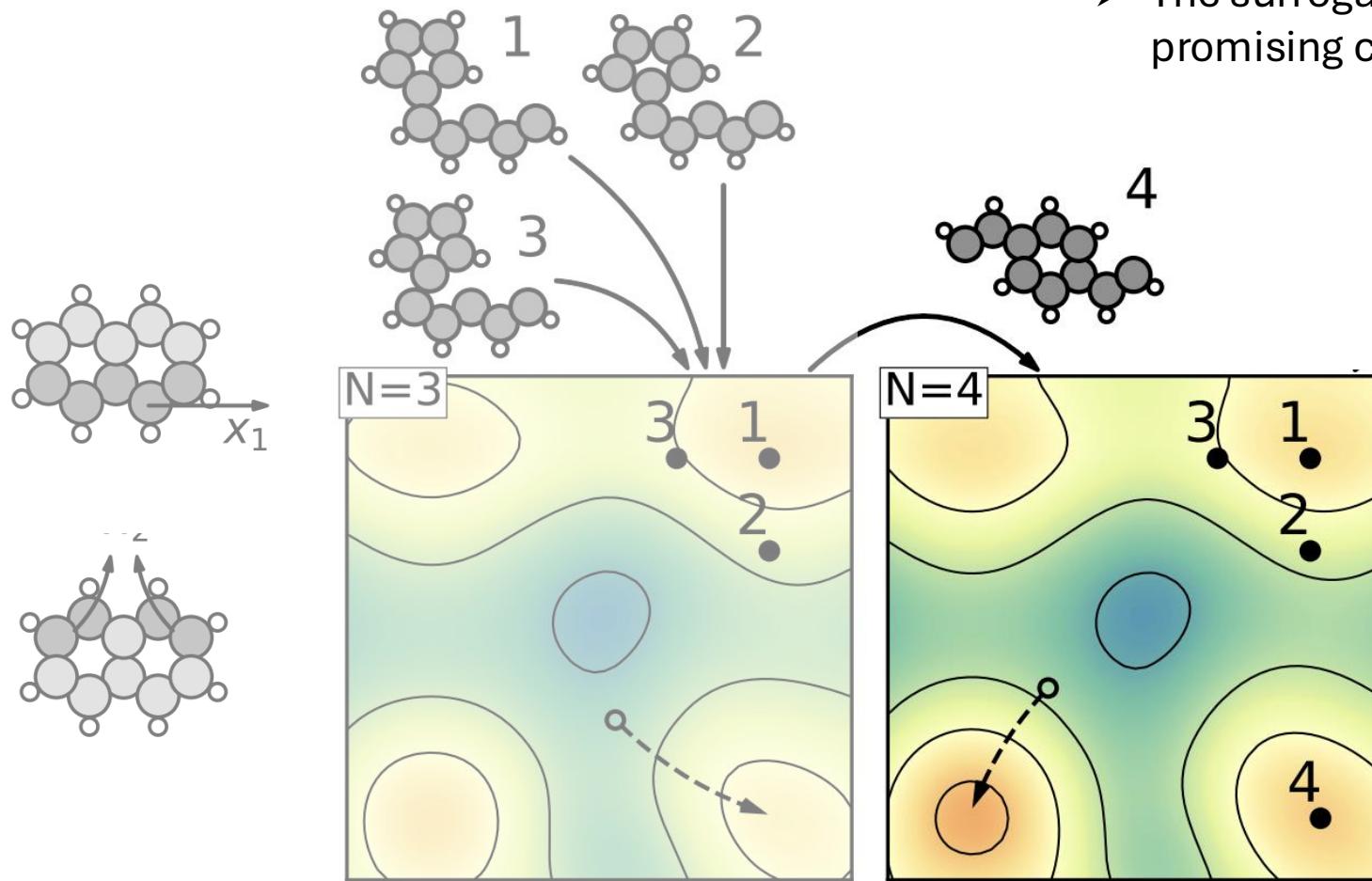


Bayesian Optimization

- The surrogate GP landscape is used to find promising configurations to calculate with DFT.

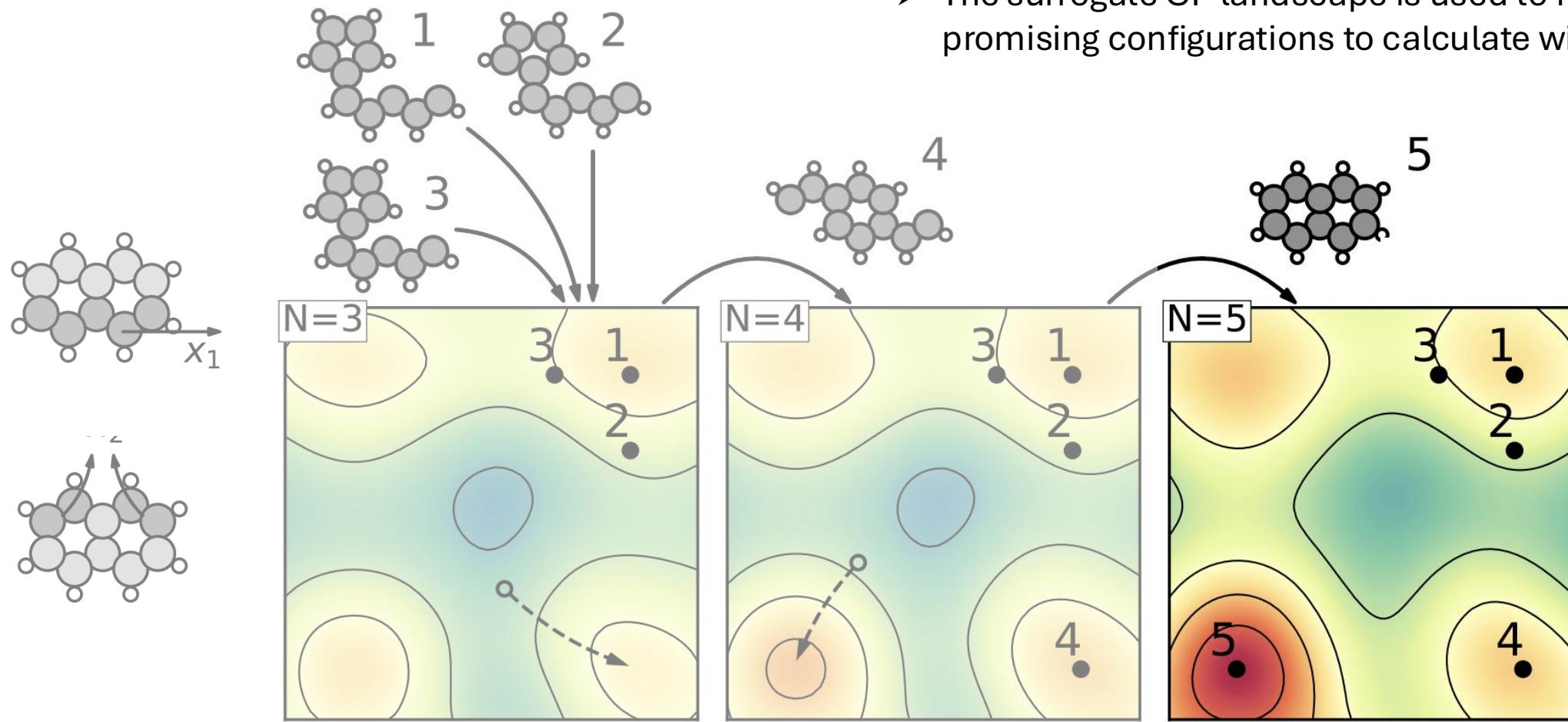
Bayesian Optimization

- The surrogate GP landscape is used to find promising configurations to calculate with DFT.

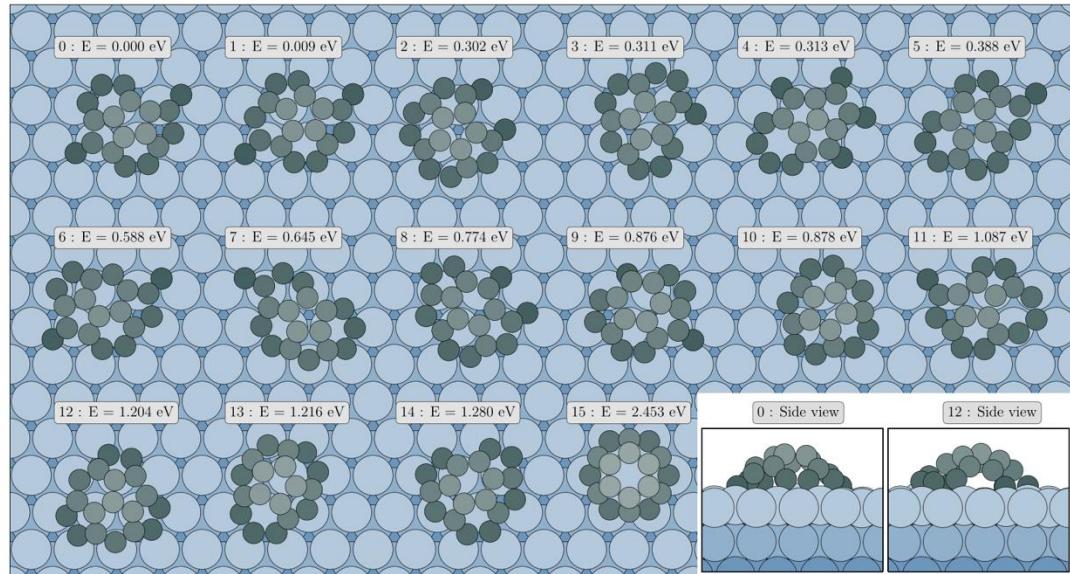


Bayesian Optimization

- The surrogate GP landscape is used to find promising configurations to calculate with DFT.

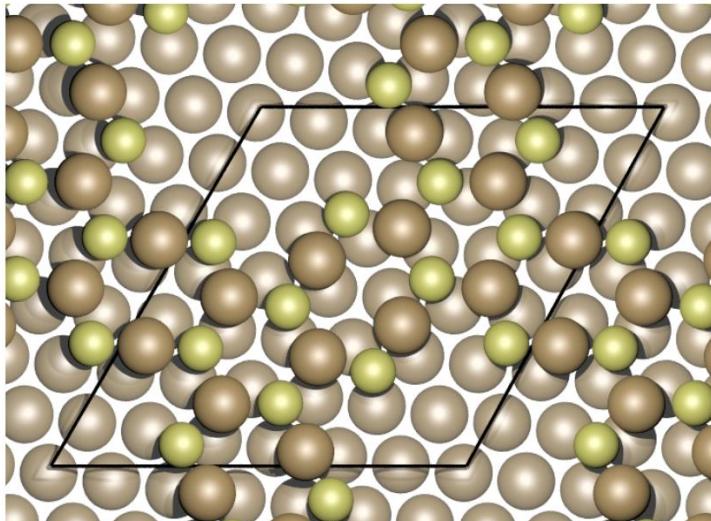


Carbon clusters on Ir(111)



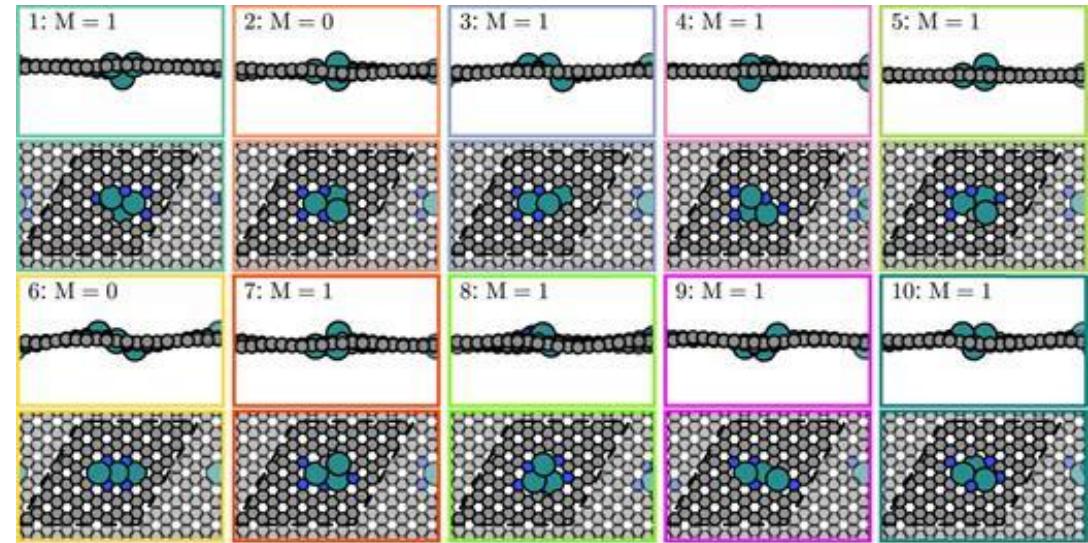
Phys. Rev. B **105**, 245404

Sulfur-induced Cu(111) reconstruction



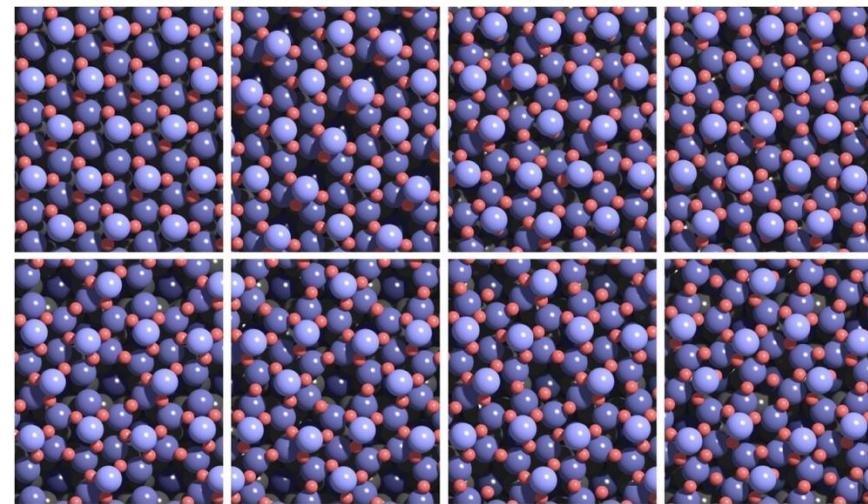
J. Chem. Phys. **160**, 174107 (2024)

Nitrogen-Ruthenium defects in graphene



J. Chem. Phys. **157**, 054701 (2022)

Tin-oxide on Pt₃Sn(111)

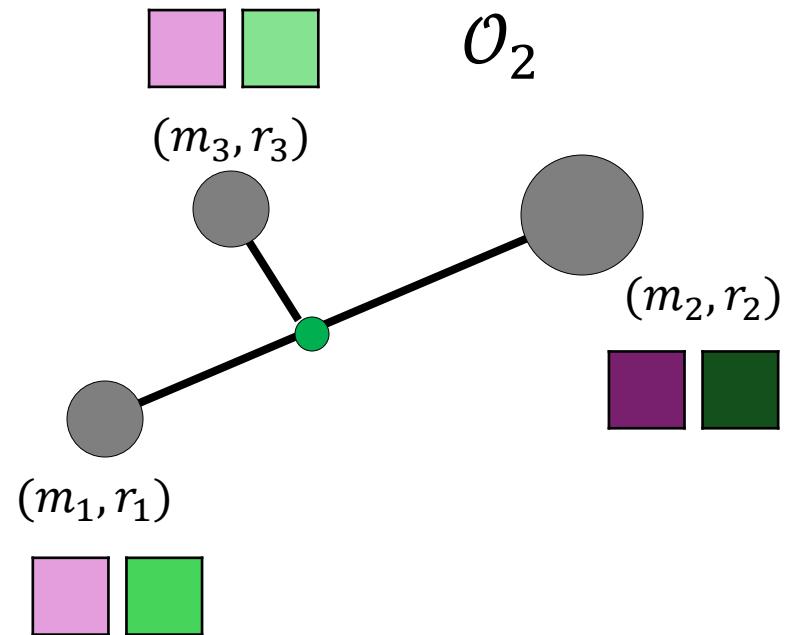
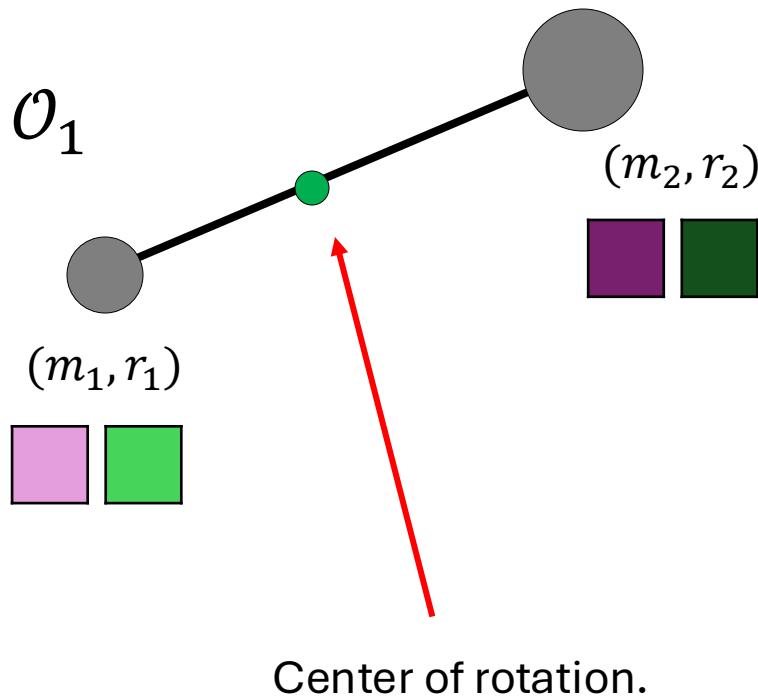


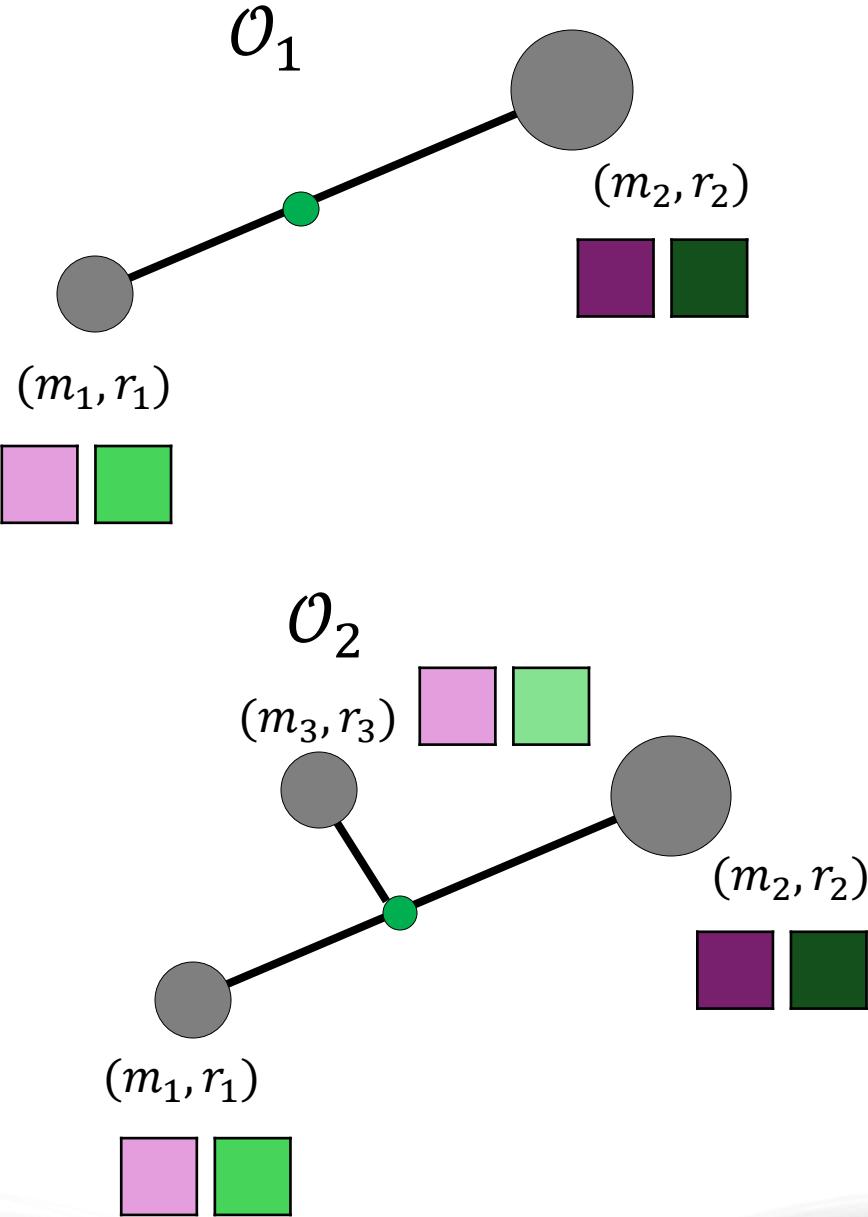
Angew. Chem. Int. Ed. **2022**, *61*, e202204244

Moment of inertia

$$I = \sum_i m_i r_i^2$$

How to calculate the covariance between these $K(\mathcal{O}_1, \mathcal{O}_2)$?





$$X_1 = \begin{array}{|c|c|} \hline \text{pink} & \text{green} \\ \hline \text{purple} & \text{dark green} \\ \hline \end{array}$$

$$X_2 = \begin{array}{|c|c|} \hline \text{pink} & \text{green} \\ \hline \text{purple} & \text{dark green} \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline \text{blue} & \text{light blue} & \text{blue} \\ \hline \text{light blue} & \text{blue} & \text{light blue} \\ \hline \end{array}$$

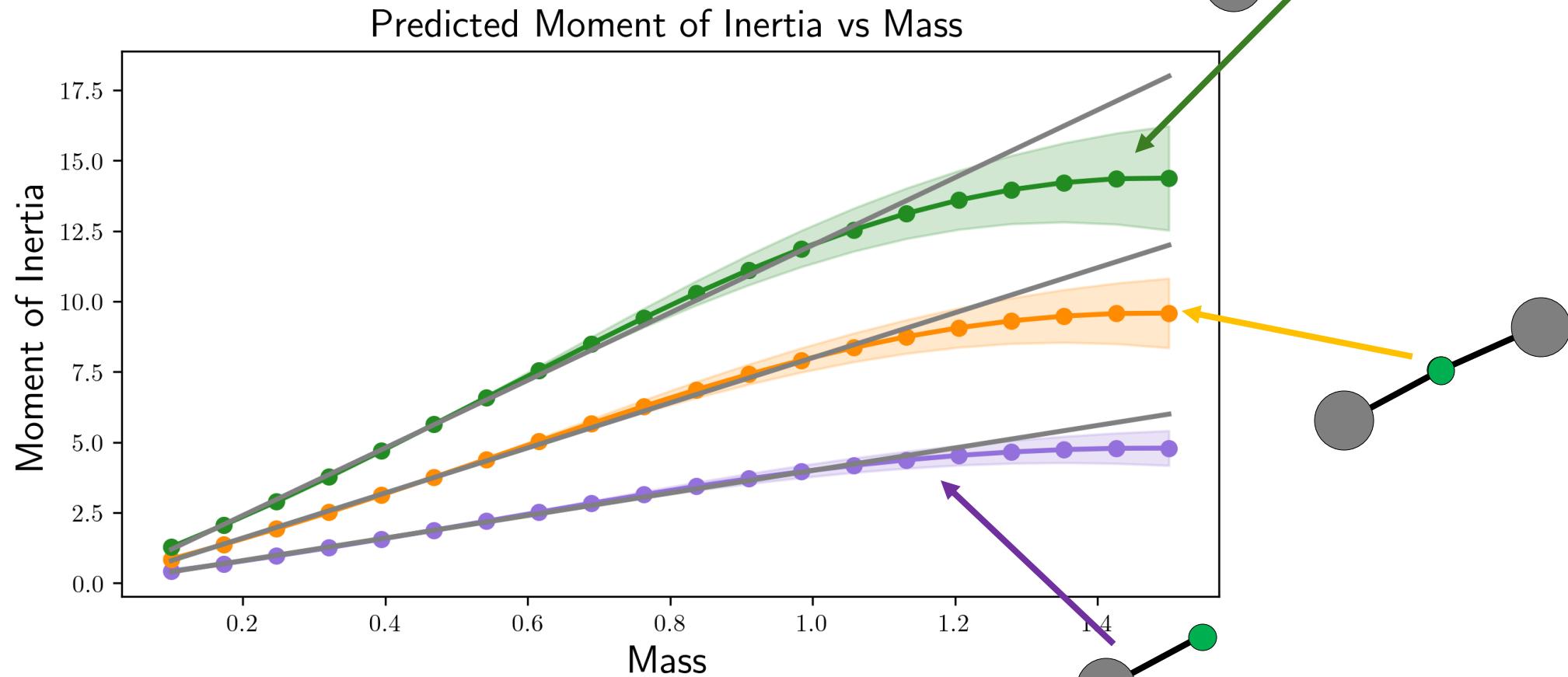
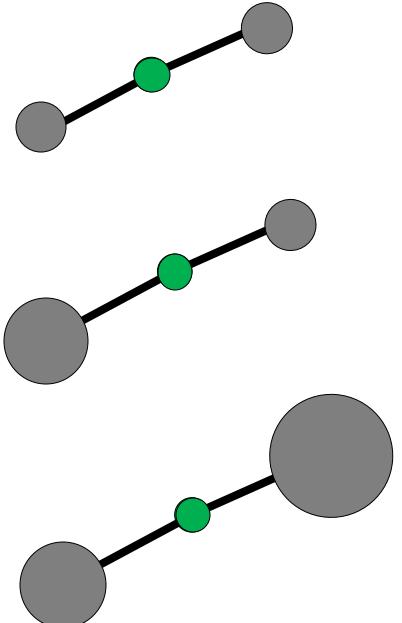
$$K(X_1, X_2) =$$

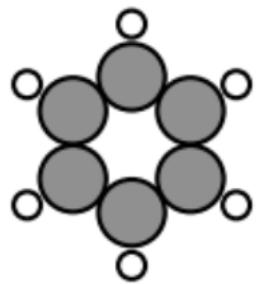
$$K(\mathcal{O}_1, \mathcal{O}_2) = \begin{array}{|c|c|} \hline 1 & 1 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline \text{blue} & \text{light blue} & \text{blue} \\ \hline \text{light blue} & \text{blue} & \text{light blue} \\ \hline \end{array} \quad \begin{array}{|c|} \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline \end{array}$$

$$= \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array}$$

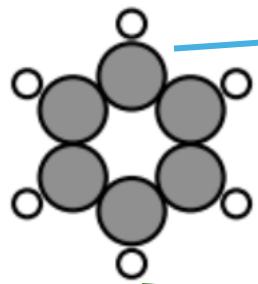
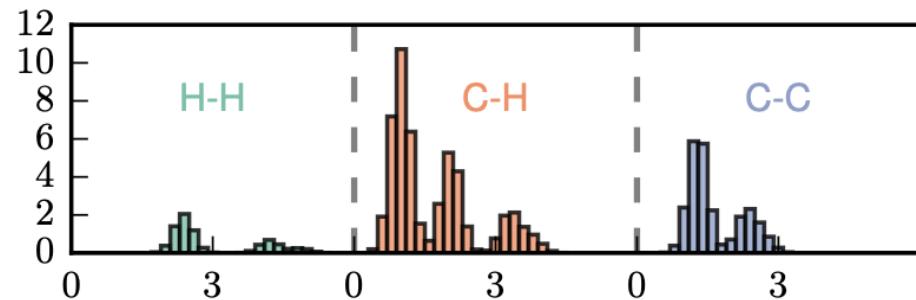
$$K(\mathcal{O}_1, \mathcal{O}_2) = L_1 K(X_1, X_2) L_2^T$$

Training examples

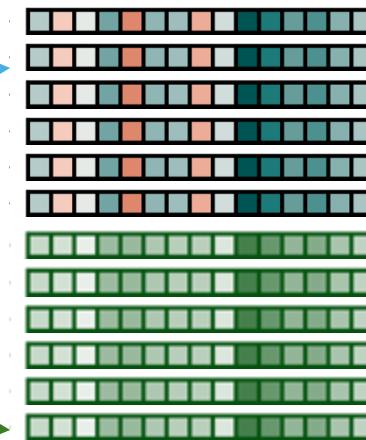




One vector

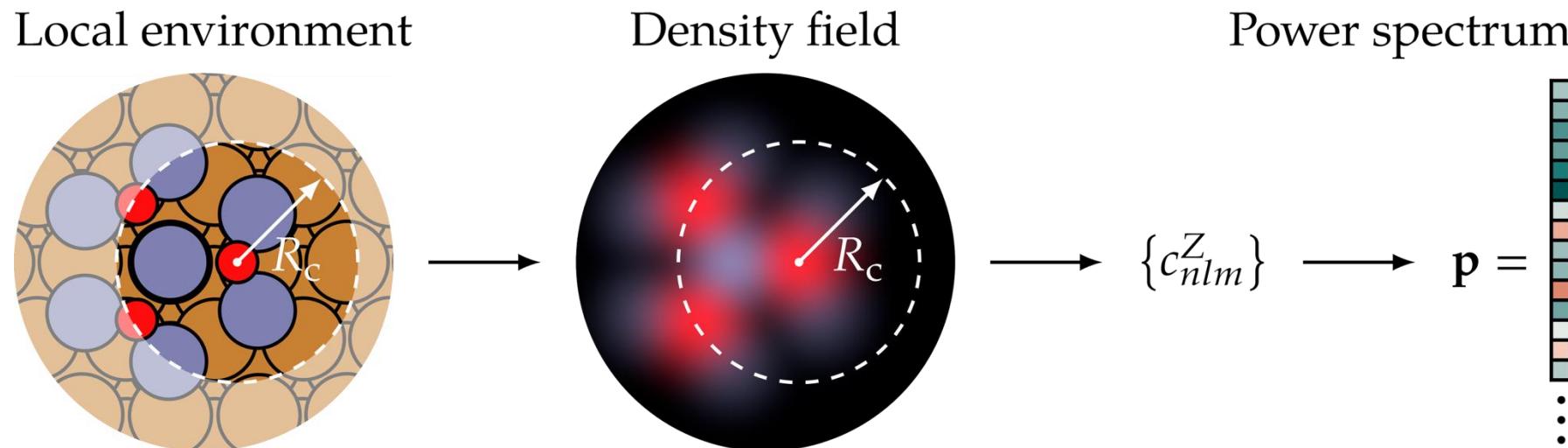


N_{atoms} vectors



Smooth Overlap of Atomic Positions (SOAP)

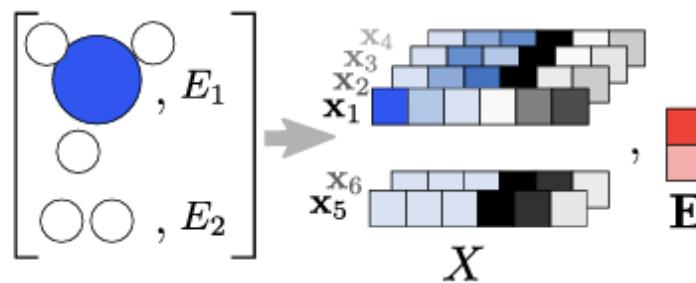
- Power spectrum of the smoothed atomic density around an atom yields an invariant descriptor.



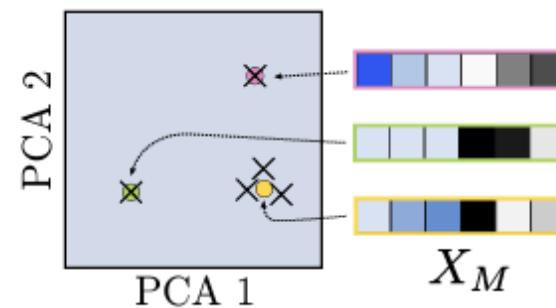
Sparse Gaussian Processes

- For large ($>1\text{-}5k$) datasets regular GPR becomes prohibitively expensive, $\mathcal{O}(n^3)$ training time and $\mathcal{O}(n)$ inference time. Sparse GPRs can be constructed with $\mathcal{O}(m^2n)$ training and $\mathcal{O}(m)$ inference where $m \ll n$.

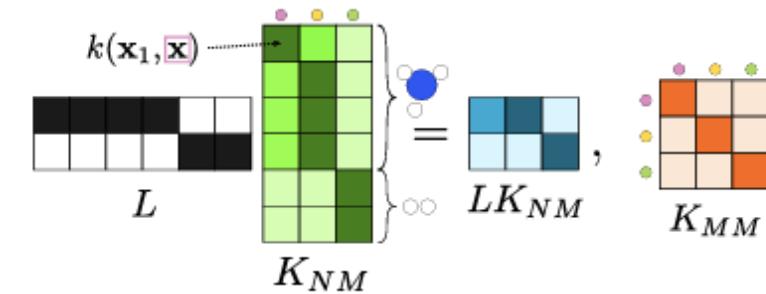
(a) Energies & atomic descriptors



(b) Basis selection



(c) Kernel



(d) Training

$$\alpha = [K_{MM} + (LK_{NM})^T \Sigma^{-1} (LK_{NM})]^{-1} (LK_{NM})^T \Sigma^{-1} \mathbf{E}$$

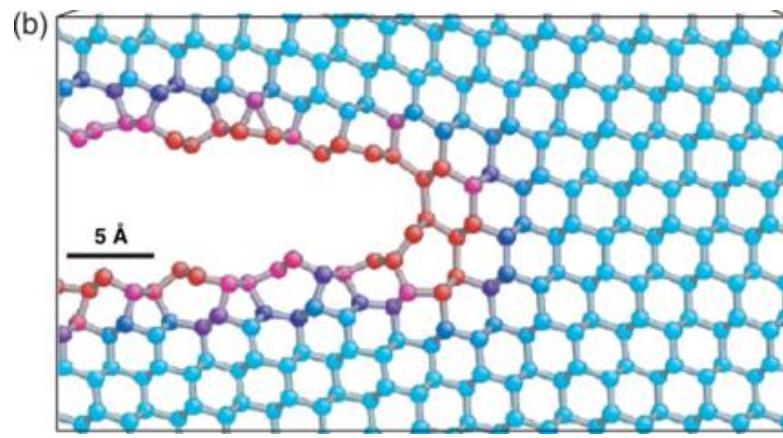
$$\begin{bmatrix} \text{Gold} \\ \text{Gold} \end{bmatrix} = \left[\begin{bmatrix} \text{Orange} & \text{Peach} & \dots \\ \text{Peach} & \text{Blue} & \dots \\ \dots & \dots & \dots \end{bmatrix} + \begin{bmatrix} \text{Grey} & \dots \\ \dots & \dots \end{bmatrix} \right]^{-1} \begin{bmatrix} \text{Blue} & \dots \\ \text{Grey} & \dots \\ \dots & \dots \end{bmatrix}^{-1} \begin{bmatrix} \text{Red} \\ \dots \end{bmatrix}$$

(e) Prediction

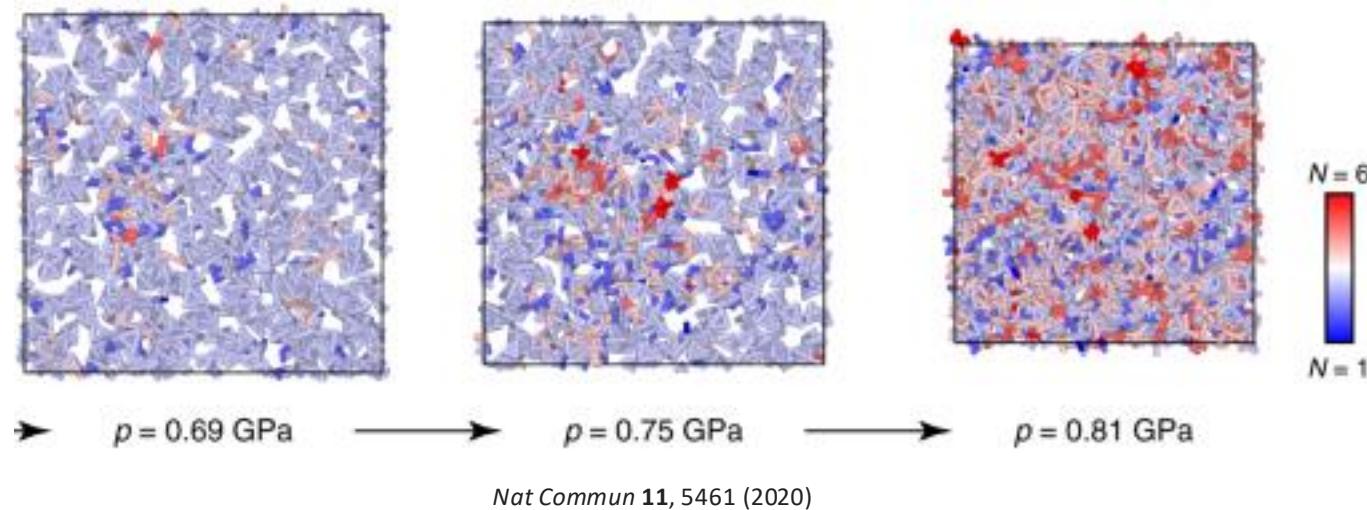
$$E_* = \sum \left[\begin{bmatrix} \text{Green} & \text{Peach} & \dots \\ \text{Peach} & \text{Blue} & \dots \\ \dots & \dots & \dots \end{bmatrix} \right]$$

The prediction process starts with a molecule with a question mark, followed by an arrow to a matrix X_*, then another arrow to the equation above.

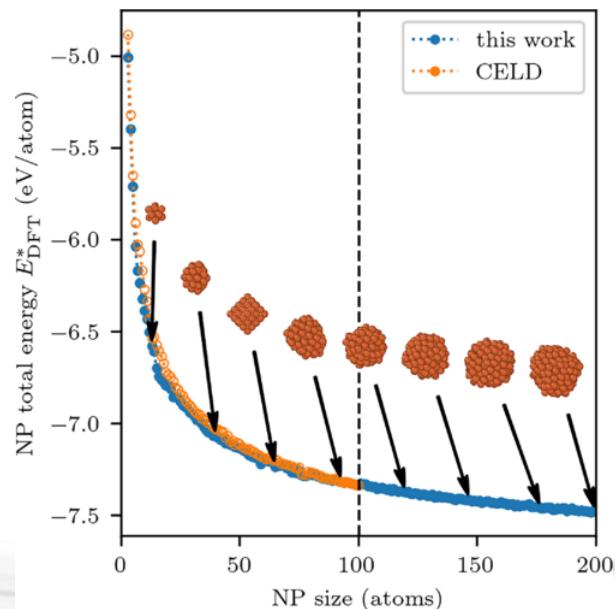
General purpose potential for Si



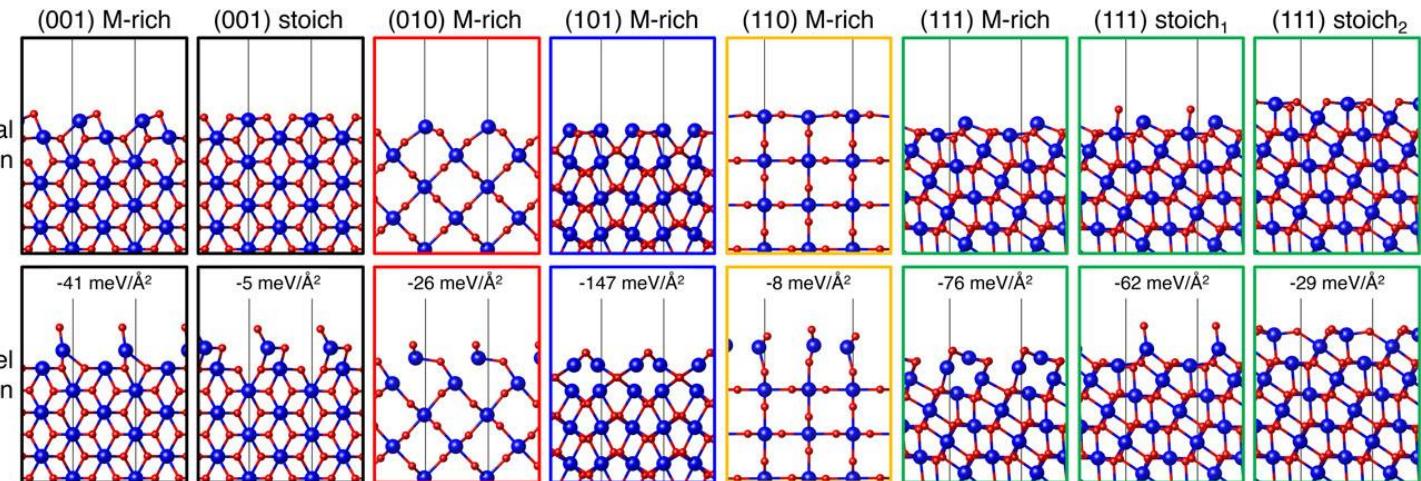
General purpose potential for phosphorus



Iron nanoparticles

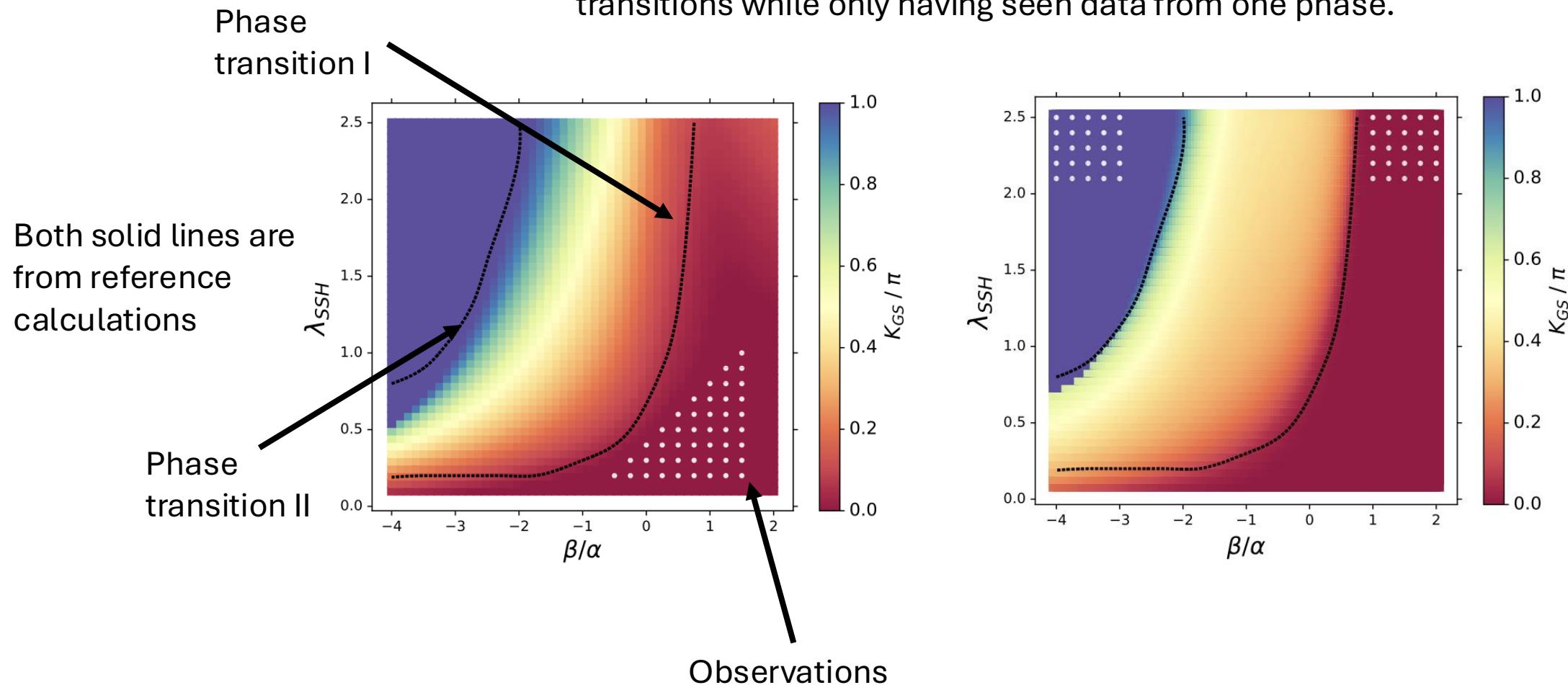


Surface structure determination of IrO_2 .



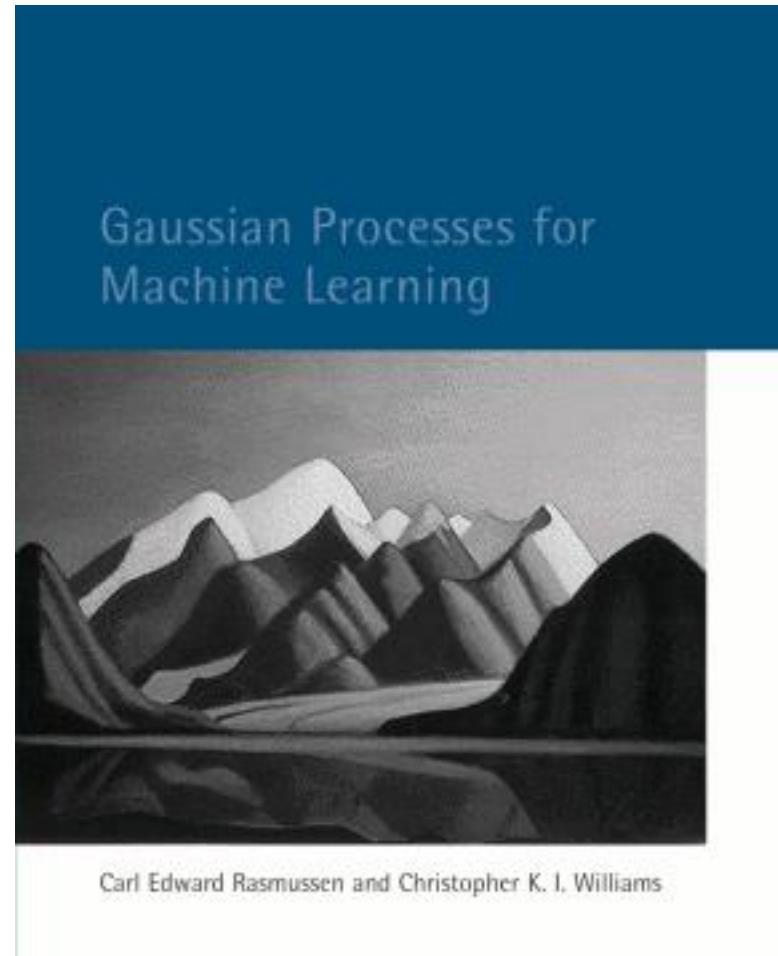
Phase transition predictions

- Optimized kernel allows GP model to predict two phase transitions while only having seen data from one phase.



Gaussian Processes for Machine Learning

- The bible of GPs for ML. PDF freely available from
<https://gaussianprocess.org/gpml/>



Packages / Libraries:

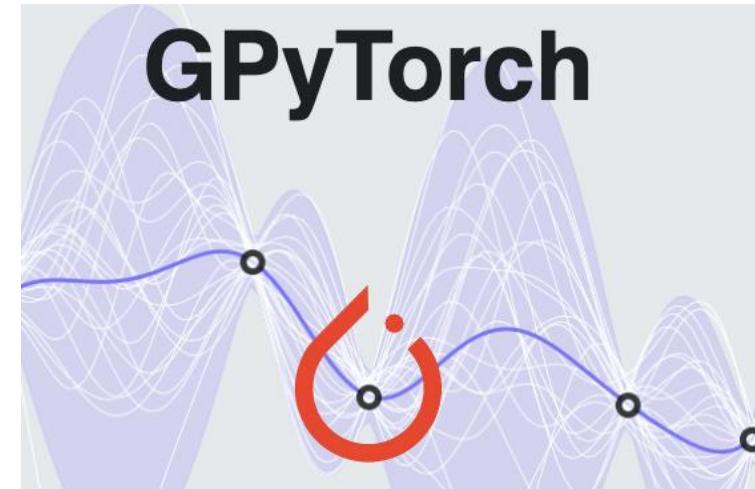
- Scikit-learn
 - Numpy implementation of GP regression.
- GPyTorch
 - Torch.
- GPJax
 - Jax.
- GPFlow
 - Tensorflow.

Among many others...

Recommendation

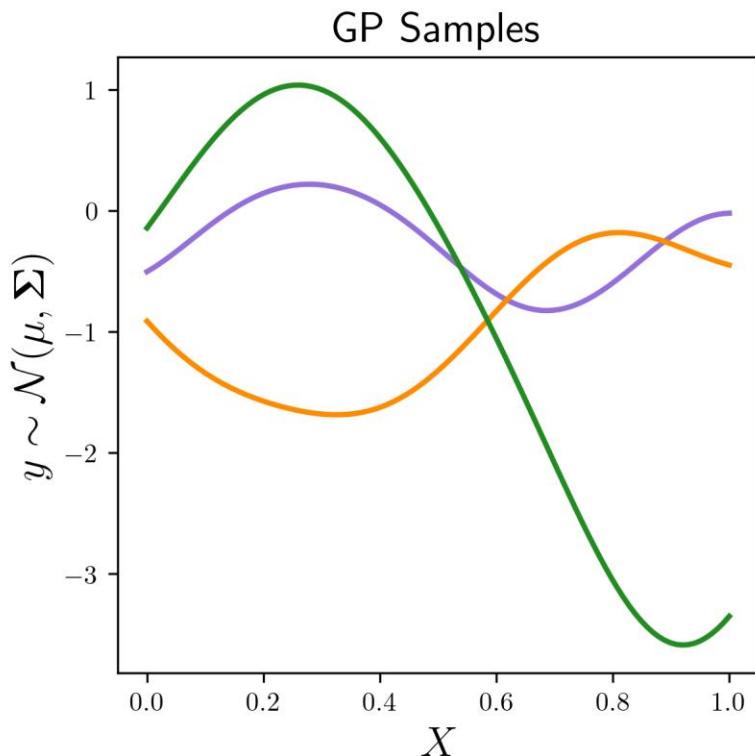
Start with your own code, read the scikit learn implementation if something is tricky (Or from my tutorial on Friday).

If you need something more feature rich or efficient try GPJax or GPyTorch depending on which backend you're more comfortable with.



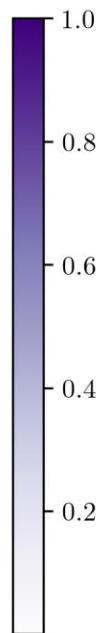
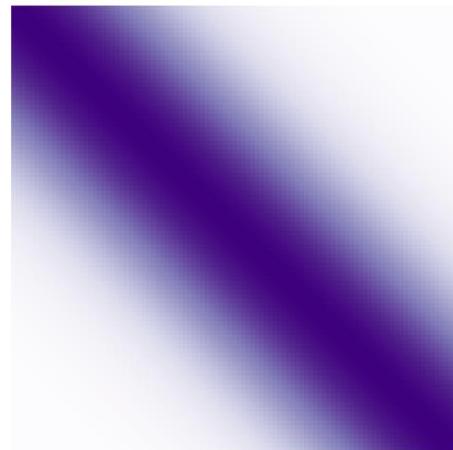
Recap

- Gaussian Processes are distributions over functions.
Uniquely defined by a mean and a covariance function



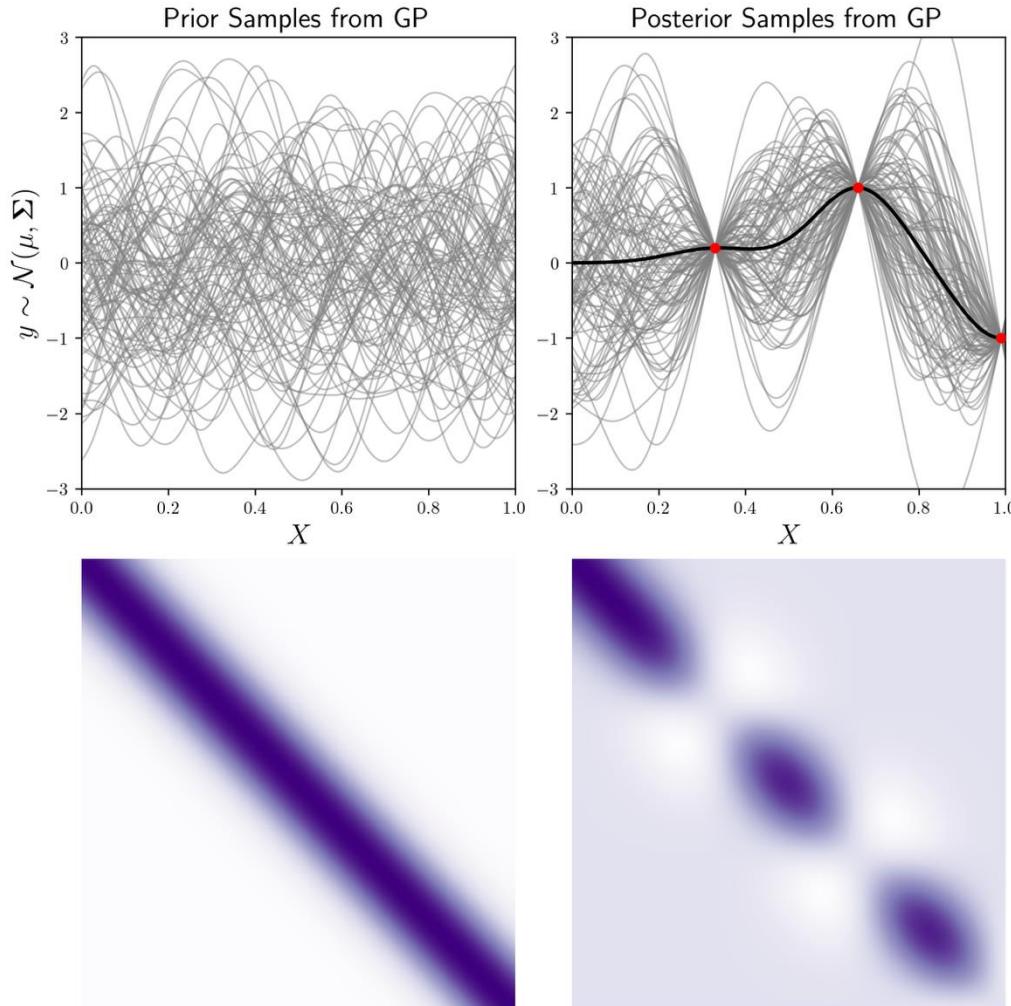
~~Covariance Matrix Σ~~

Covariance function



Recap

- Conditioning tightens the distribution such that only relevant functions are kept.



Acknowledgements



**Funded by
the European Union**



- This work is based upon work from COST Action CA21126 - Carbon molecular nanostructures in space (NanoSpace), supported by COST (European Cooperation in Science & Technology).
- COST (European Cooperation in Science & Technology) is a funding agency for research and innovation networks. Our Actions help connect research initiatives across Europe and enable scientists to grow their ideas by sharing them with their peers. This boosts their research, career and innovation.
- www.cost.eu
- <https://www.cost.eu/actions/CA21126/>
- <https://research.iac.es/proyecto/nanospace/>