

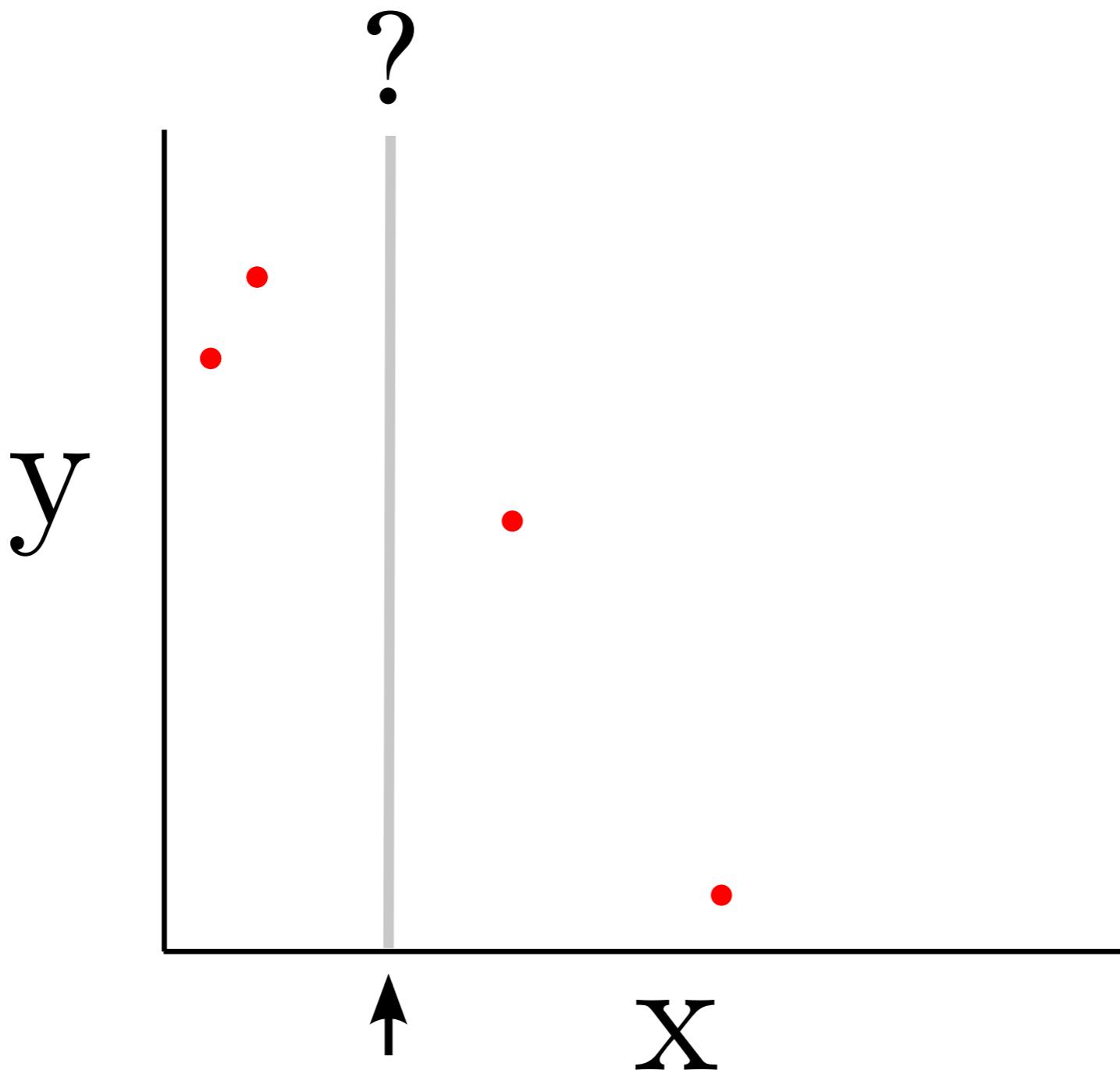
# DS-GA 3001.008 Modeling time series data

## L11. GPs

Instructor: Cristina Savin  
NYU, CNS & CDS

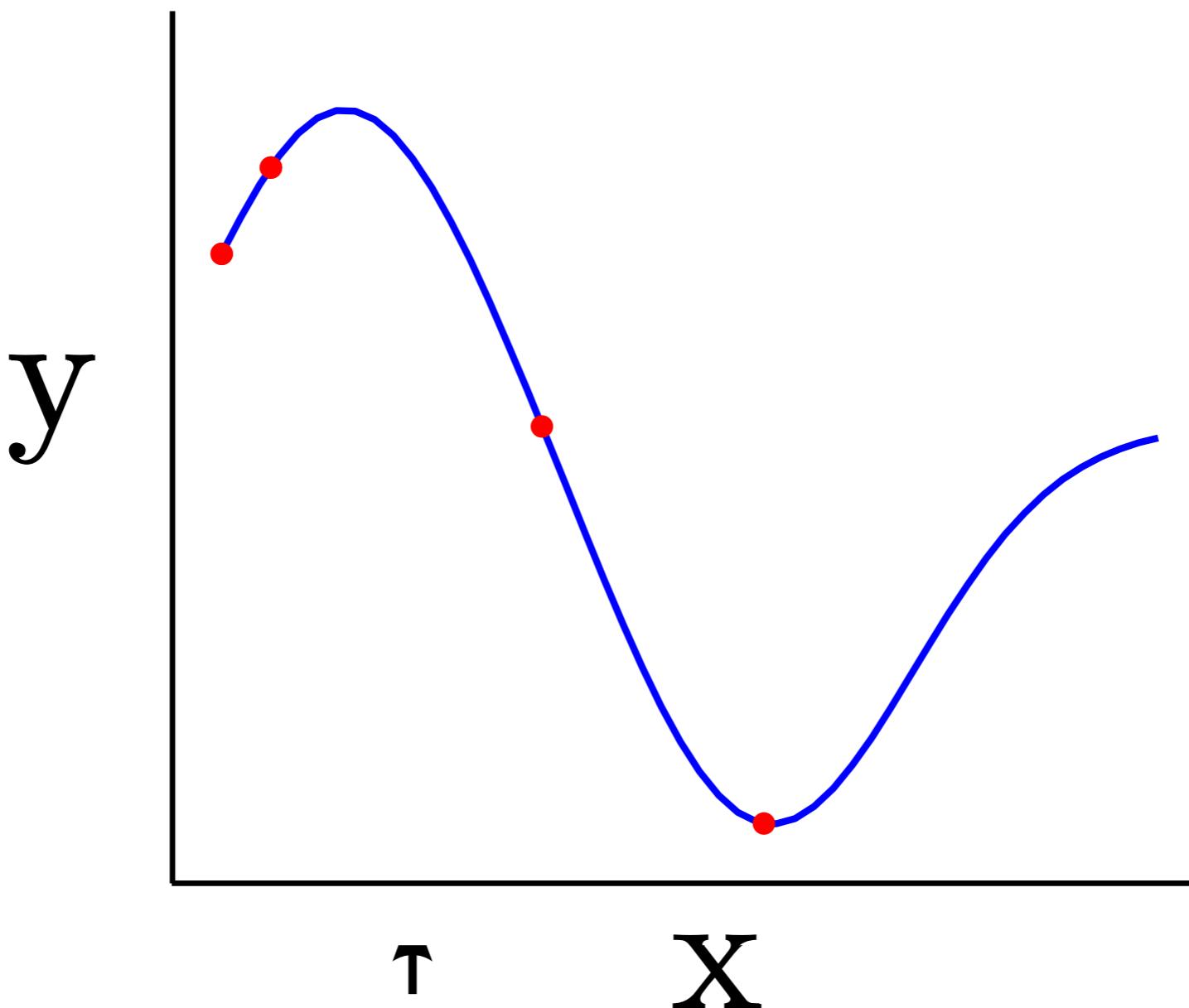
## Quick recap

### 1. Motivation: nonlinear regression



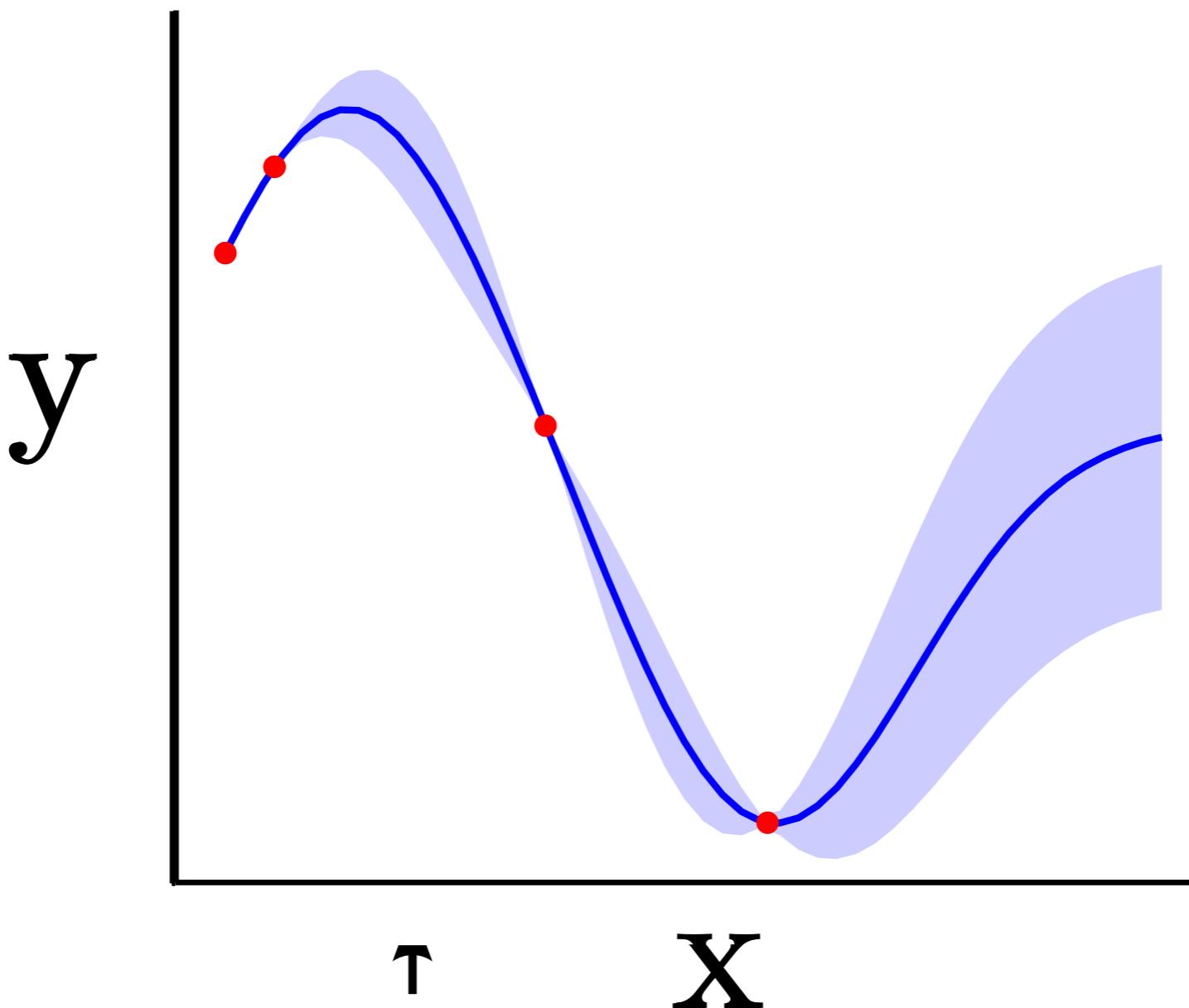
## Quick recap

### 1. Motivation: nonlinear regression



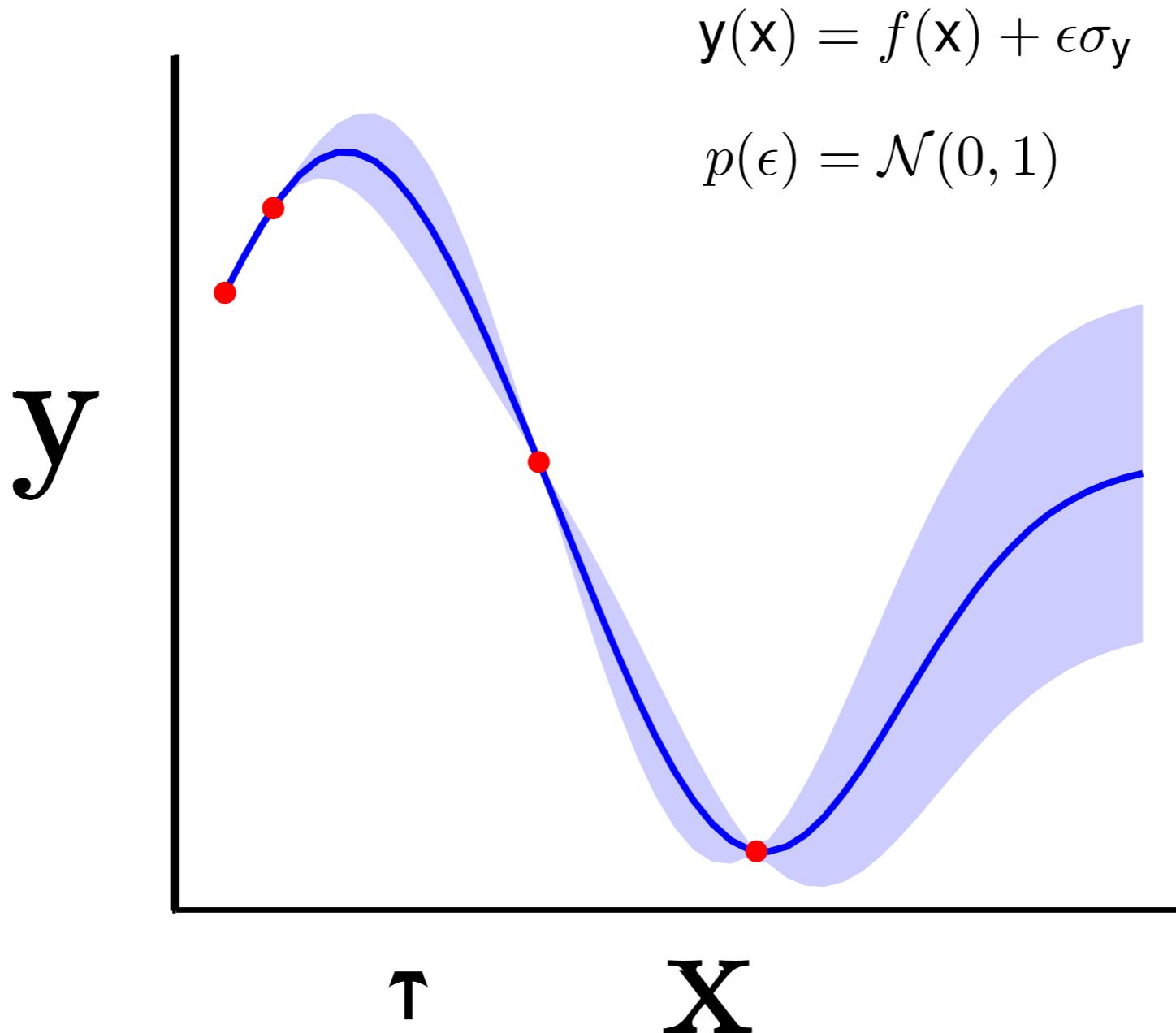
## Quick recap

### 1. Motivation: nonlinear regression



## Quick recap

### 1. Motivation: nonlinear regression

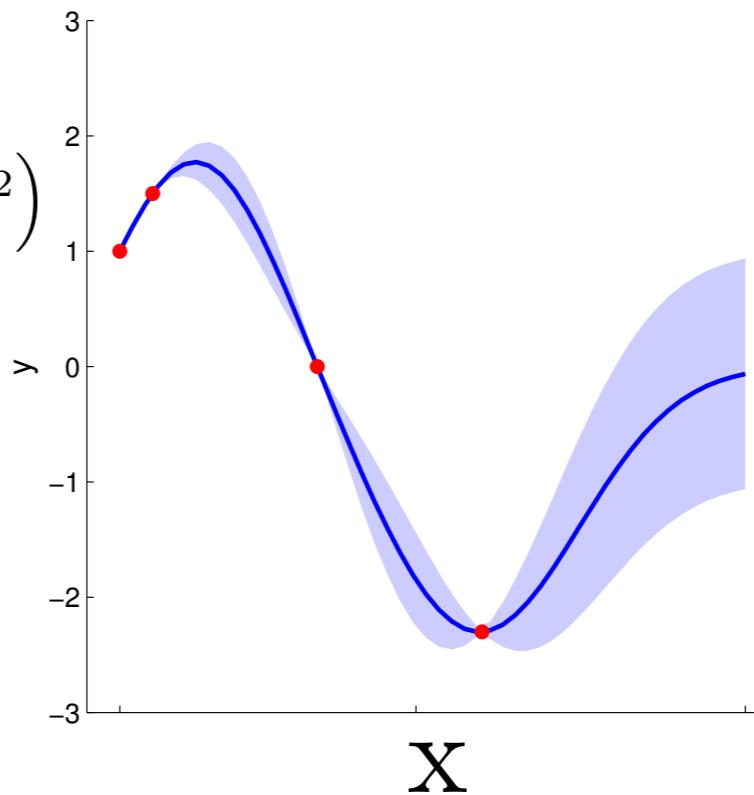
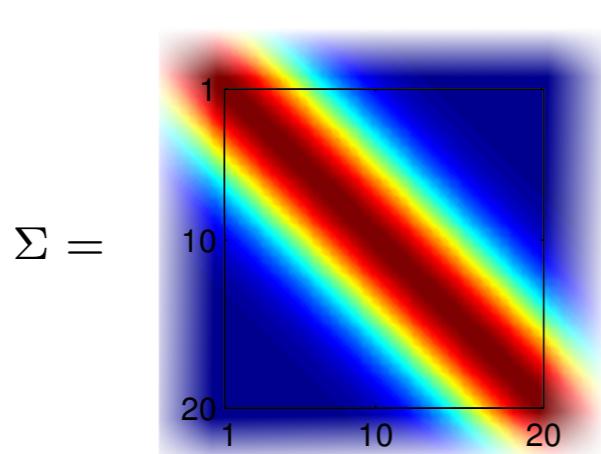


## Quick recap

### 2.prior over functions: GP

$$\Sigma(\mathbf{x}_1, \mathbf{x}_2) = K(\mathbf{x}_1, \mathbf{x}_2) + I\sigma_y^2$$

$$K(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(\mathbf{x}_1 - \mathbf{x}_2)^2\right)$$



**GP:** generalization of multivariate gaussian

Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

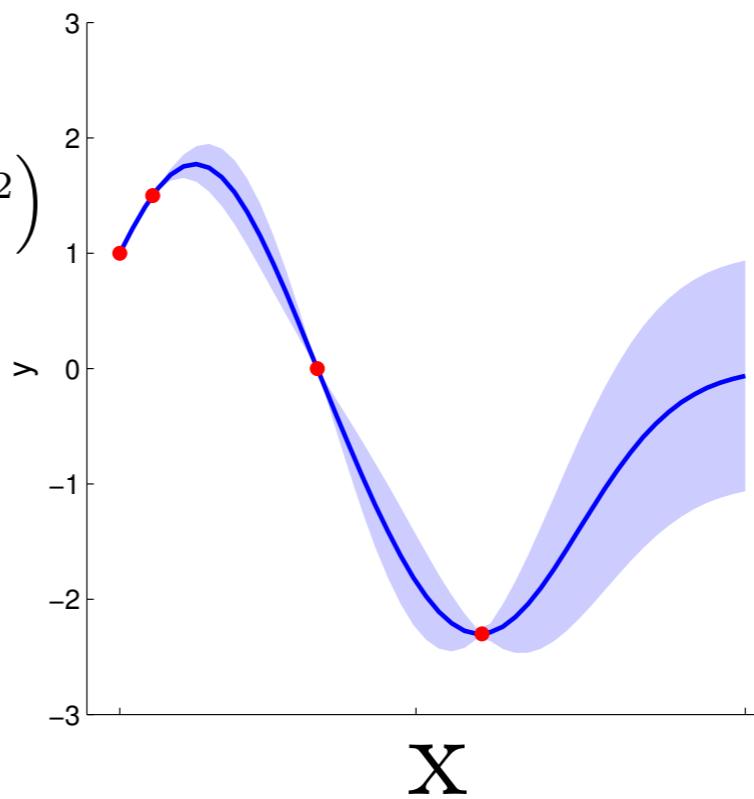
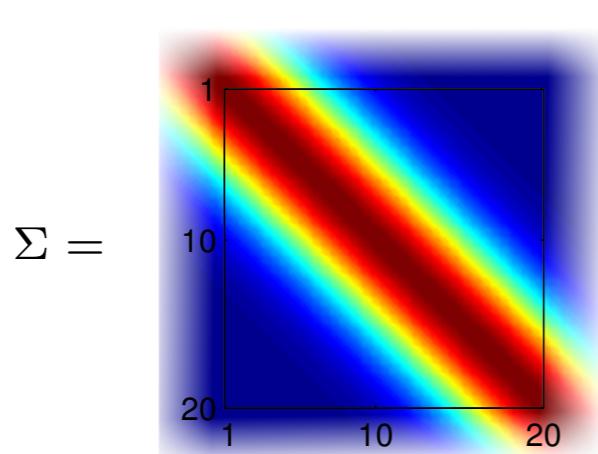
$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')) \quad \text{mean+covariance functions}$$

## Quick recap

### 2.prior over functions: GP

$$\Sigma(\mathbf{x}_1, \mathbf{x}_2) = K(\mathbf{x}_1, \mathbf{x}_2) + I\sigma_y^2$$

$$K(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(\mathbf{x}_1 - \mathbf{x}_2)^2\right)$$



**GP:** generalization of multivariate gaussian

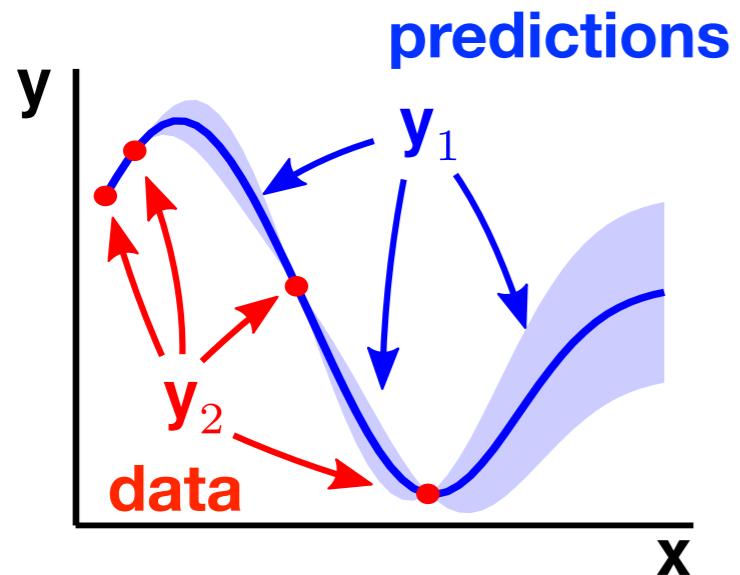
Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')) \quad \text{mean+covariance functions}$$

Any function has nonzero probability, preference for certain structure

## Quick recap

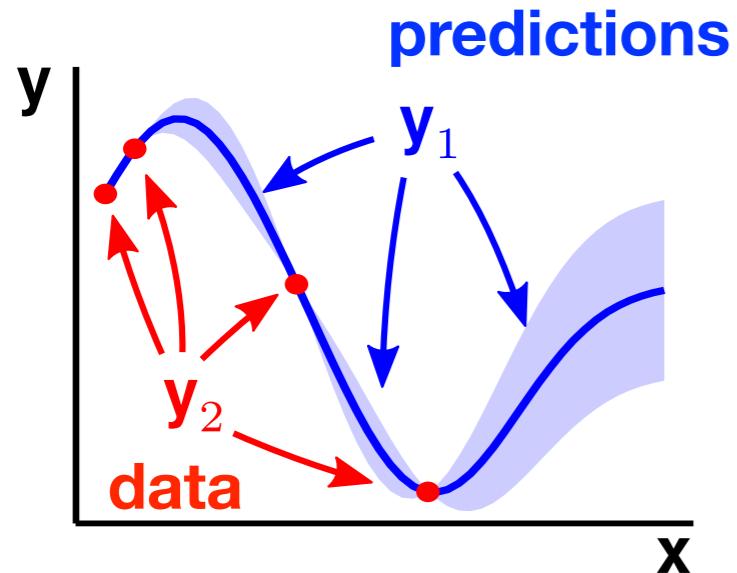
### 3. GP inference



$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)}$$

## Quick recap

### 3. GP inference



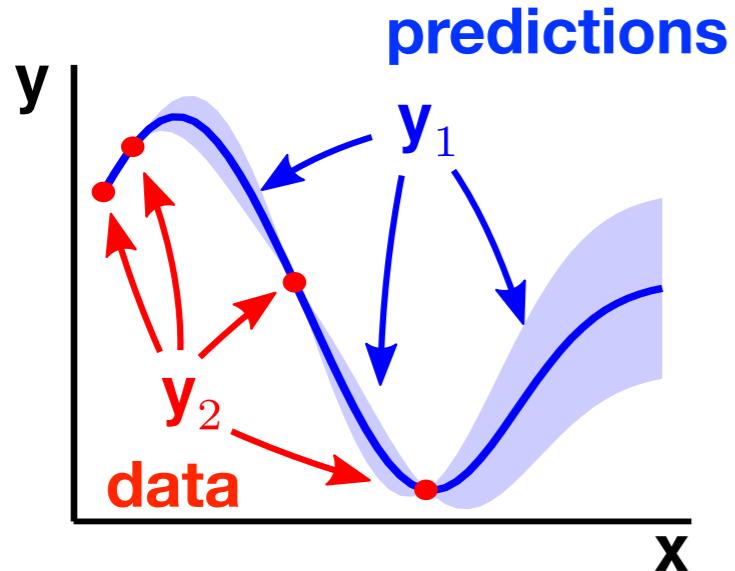
$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)}$$

**Jointly gaussian:**

$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left( \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \right)$$

## Quick recap

### 3. GP inference



$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)}$$

**Jointly gaussian:**

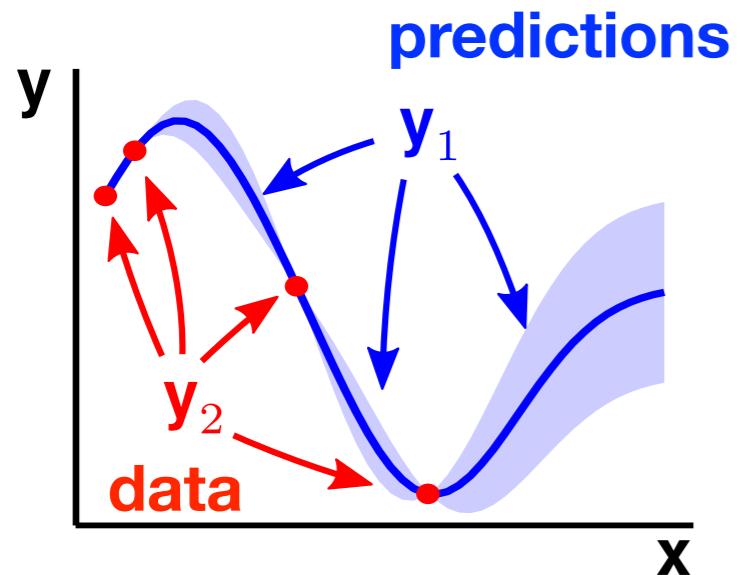
$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left( \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \right)$$

↓

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\mathbf{a} + \mathbf{B}\mathbf{C}^{-1}(\mathbf{y}_2 - \mathbf{b}), \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T)$$

## Quick recap

### 3. GP inference



$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)}$$

**Jointly gaussian:**

$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left( \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \right)$$

↓

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\mathbf{a} + \mathbf{B}\mathbf{C}^{-1}(\mathbf{y}_2 - \mathbf{b}), \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T)$$

**predicted mean**

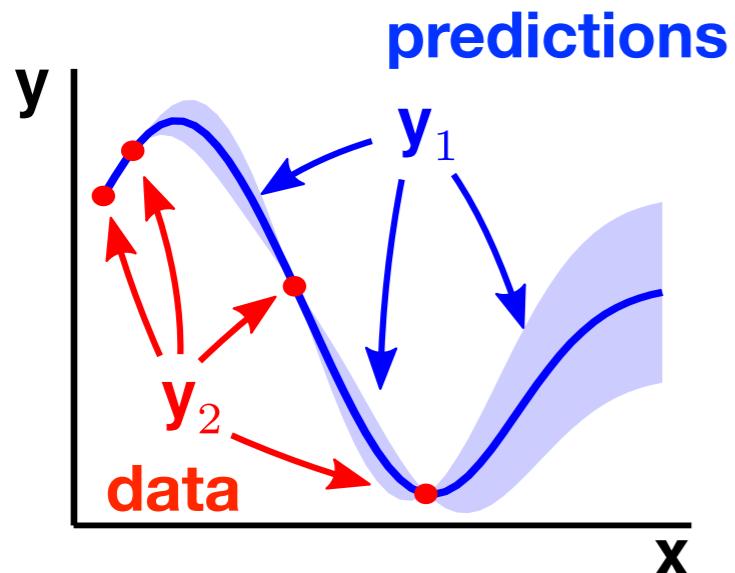
$$\begin{aligned} \mu_{\mathbf{y}_1 | \mathbf{y}_2} &= \mathbf{a} + \mathbf{B}\mathbf{C}^{-1}(\mathbf{y}_2 - \mathbf{b}) \\ &= \mathbf{B}\mathbf{C}^{-1}\mathbf{y}_2 \\ &= \mathbf{W}\mathbf{y}_2 \end{aligned}$$

**predicted covariance**

$$\Sigma_{\mathbf{y}_1 | \mathbf{y}_2} = \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T$$

## Quick recap

### 3. GP inference



$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)}$$

**Jointly gaussian:**

$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left( \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \right)$$



$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\mathbf{a} + \mathbf{BC}^{-1}(\mathbf{y}_2 - \mathbf{b}), \mathbf{A} - \mathbf{BC}^{-1}\mathbf{B}^T)$$

**predicted mean**

$$\begin{aligned} \mu_{\mathbf{y}_1 | \mathbf{y}_2} &= \mathbf{a} + \mathbf{BC}^{-1}(\mathbf{y}_2 - \mathbf{b}) \\ &= \mathbf{BC}^{-1}\mathbf{y}_2 \\ &= \mathbf{W}\mathbf{y}_2 \end{aligned}$$

**predicted covariance**

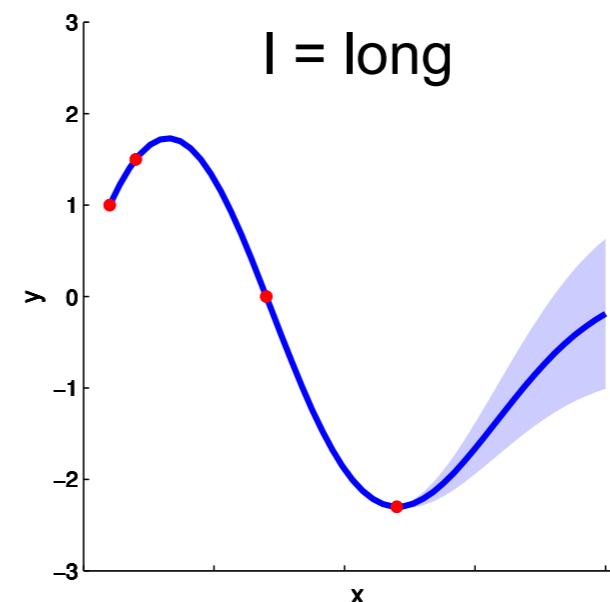
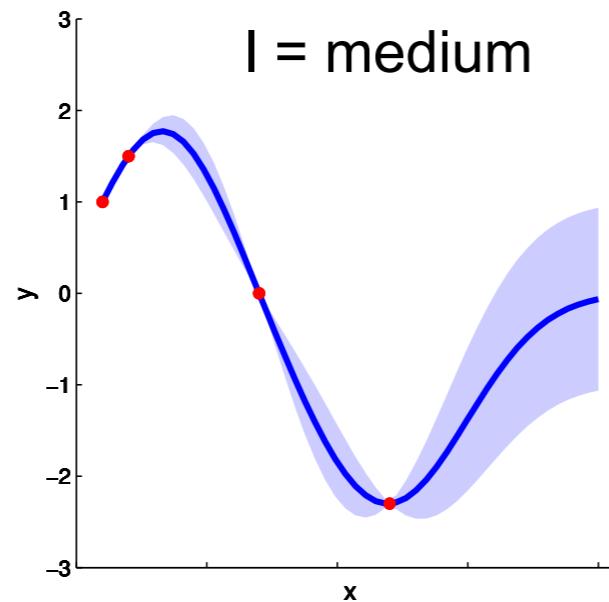
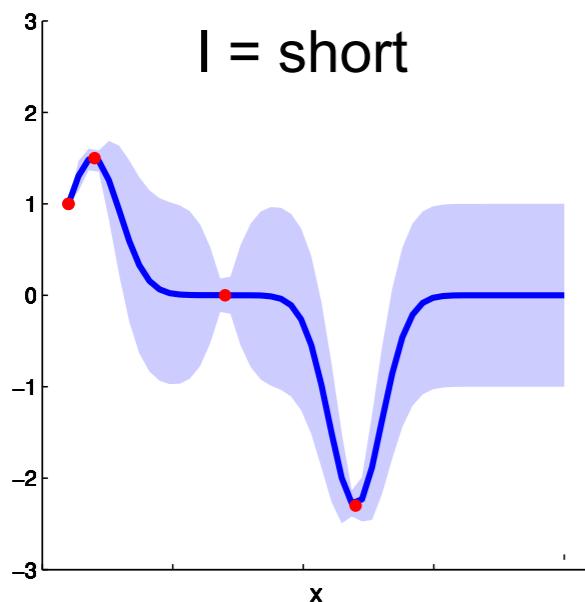
$$\Sigma_{\mathbf{y}_1 | \mathbf{y}_2} = \mathbf{A} - \mathbf{BC}^{-1}\mathbf{B}^T$$

**Computational bottleneck!!!**

## Quick recap

### 4. GP hyperparameters

$$K(x_1, x_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(x_1 - x_2)^2\right)$$

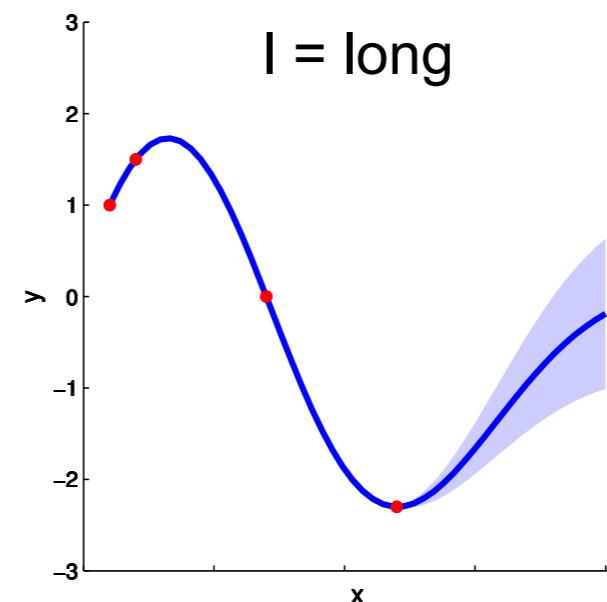
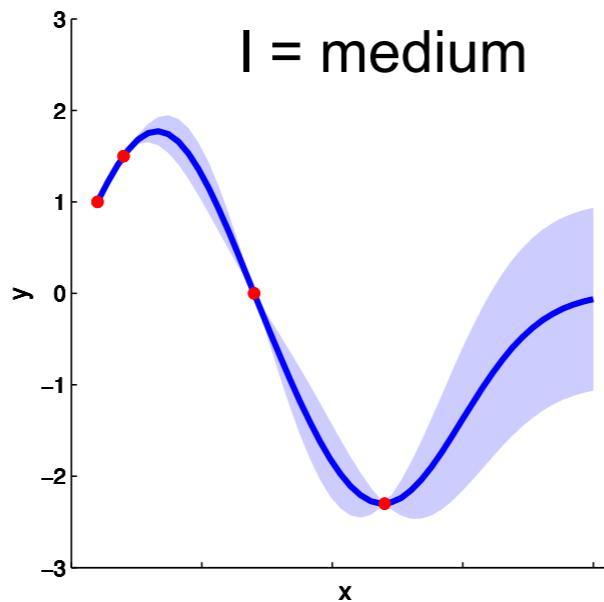
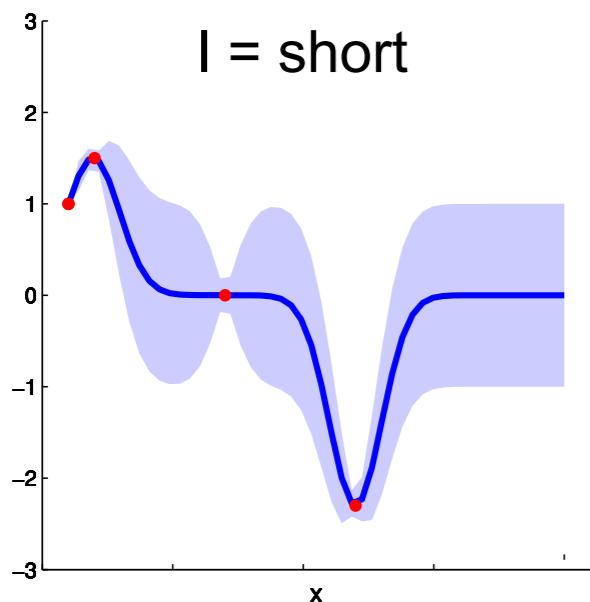


Hyperparameters  
significantly  
influence outcome

## Quick recap

### 4. GP hyperparameters

$$K(x_1, x_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(x_1 - x_2)^2\right)$$



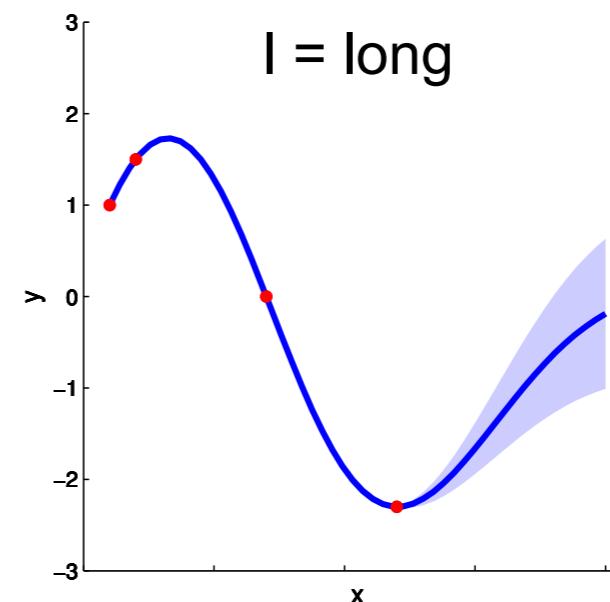
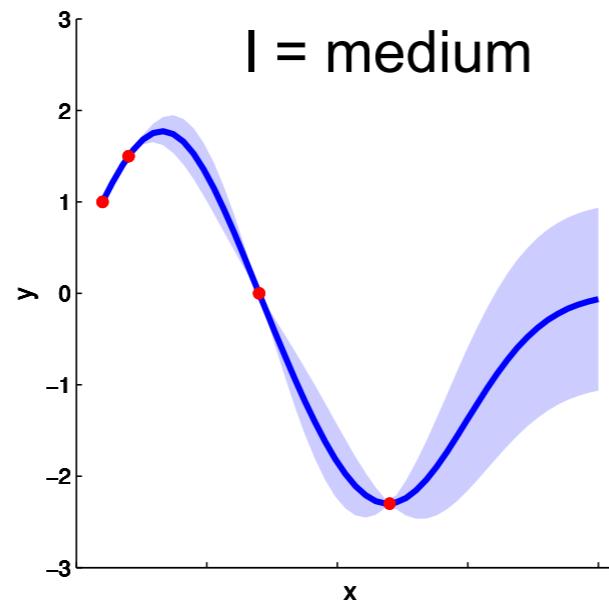
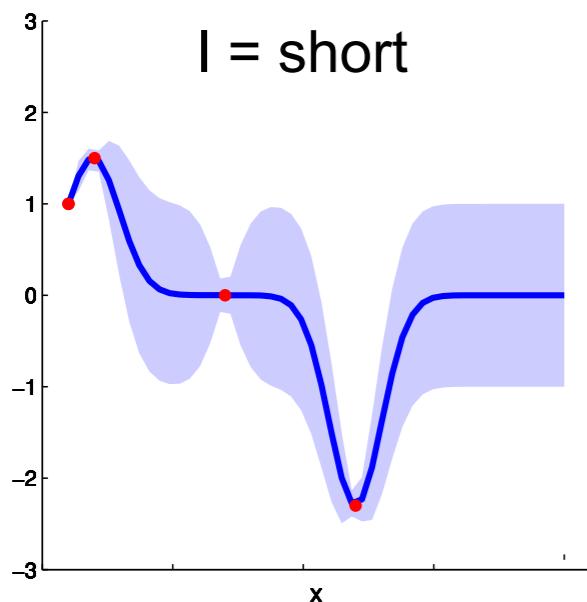
Hyperparameters significantly influence outcome

Learn values from data!

## Quick recap

### 4. GP hyperparameters

$$K(x_1, x_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(x_1 - x_2)^2\right)$$



Hyperparameters significantly influence outcome

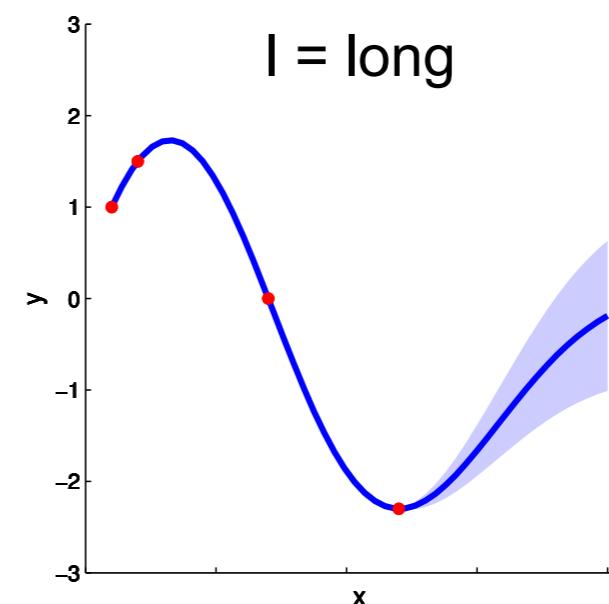
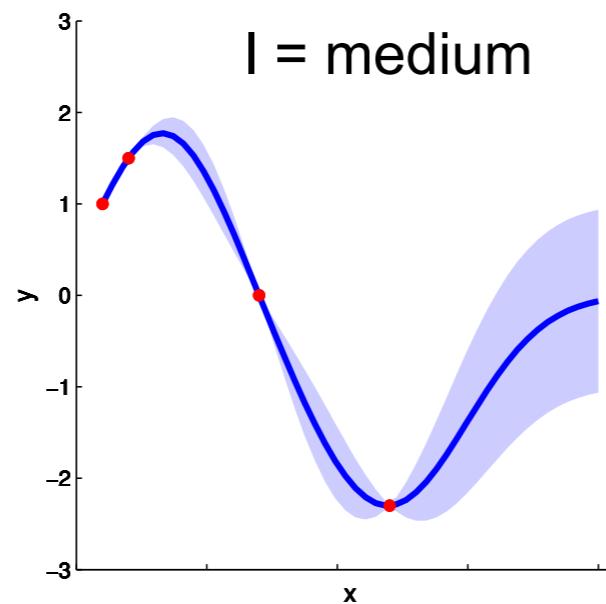
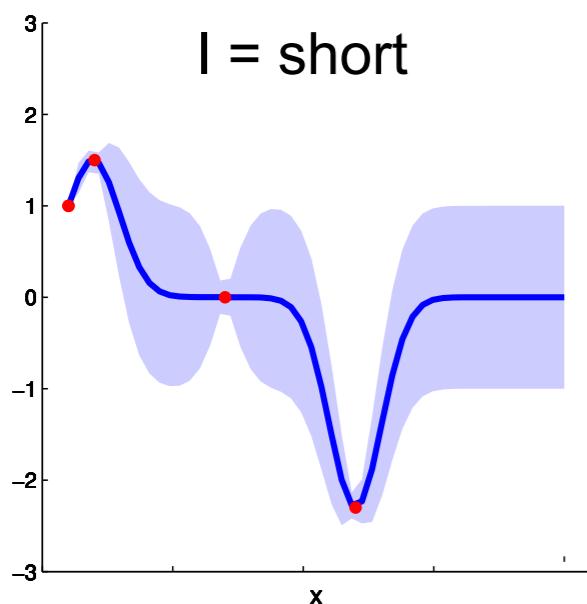
Learn values from data!

A. Maximum likelihood fit:

## Quick recap

### 4. GP hyperparameters

$$K(x_1, x_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(x_1 - x_2)^2\right)$$



Hyperparameters significantly influence outcome

Learn values from data!

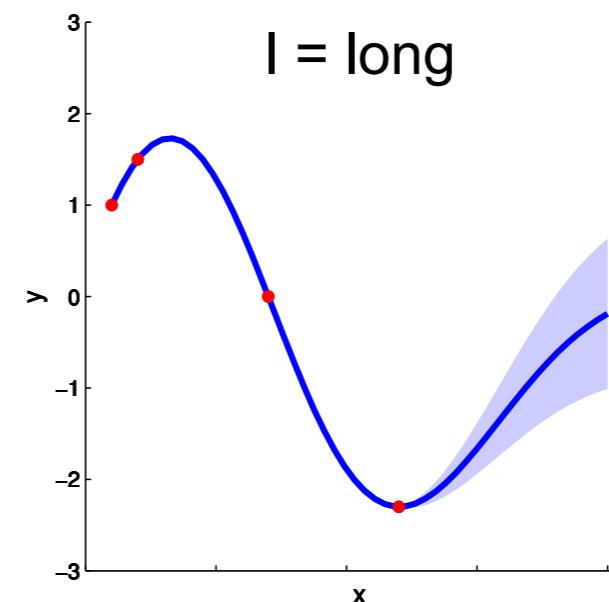
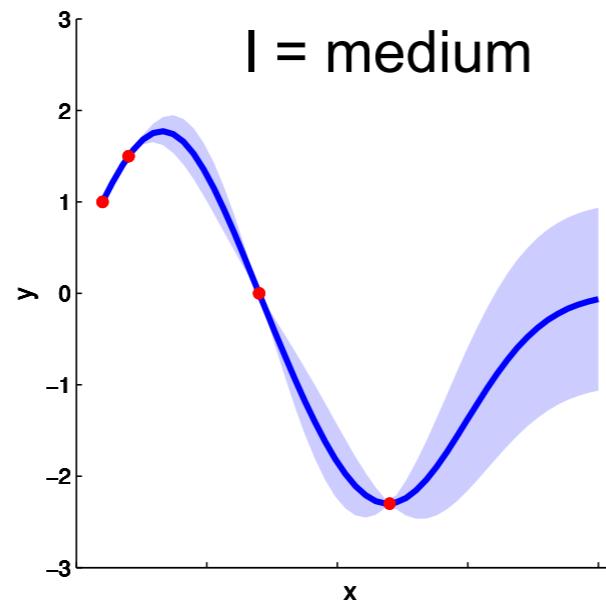
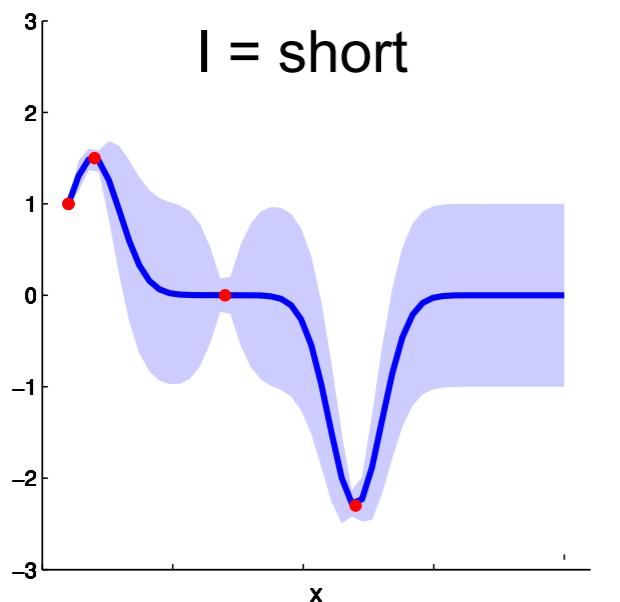
A. Maximum likelihood fit:

$$\operatorname{argmax}_{\theta} P(\mathbf{y} | \theta)$$

## Quick recap

### 4. GP hyperparameters

$$K(x_1, x_2) = \sigma^2 \exp\left(-\frac{1}{2l^2}(x_1 - x_2)^2\right)$$



Hyperparameters significantly influence outcome

Learn values from data!

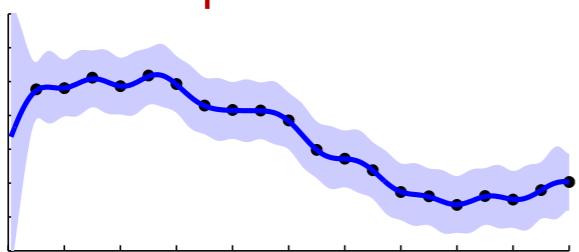
**A. Maximum likelihood fit:**

$$\operatorname{argmax}_{\theta} P(\mathbf{y} | \theta)$$

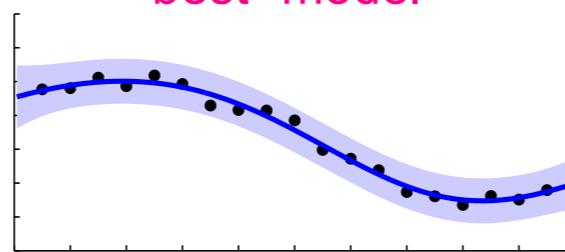
**B. Bayesian:**

$$p(\theta | \mathbf{y}_{1:N}) = \frac{p(\mathbf{y}_{1:N} | \theta) p(\theta)}{p(\mathbf{y}_{1:N})}$$

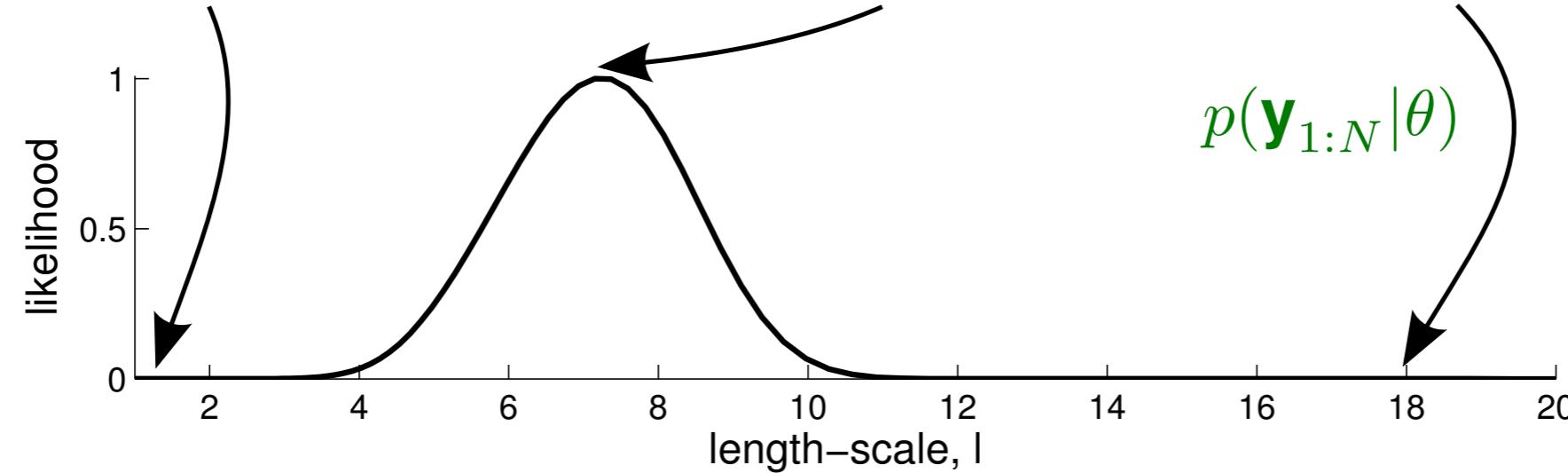
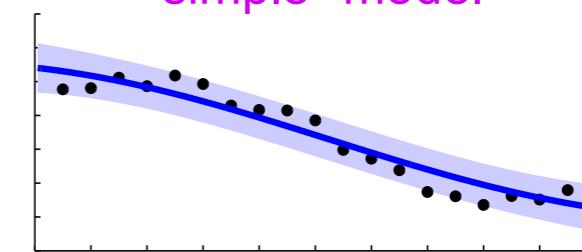
fits every training point  
"complex" model



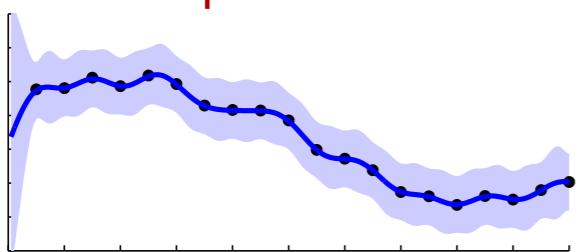
"best" model



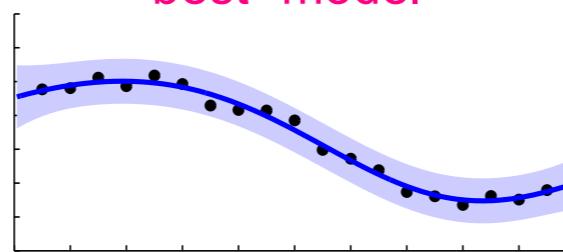
straight line-like  
"simple" model



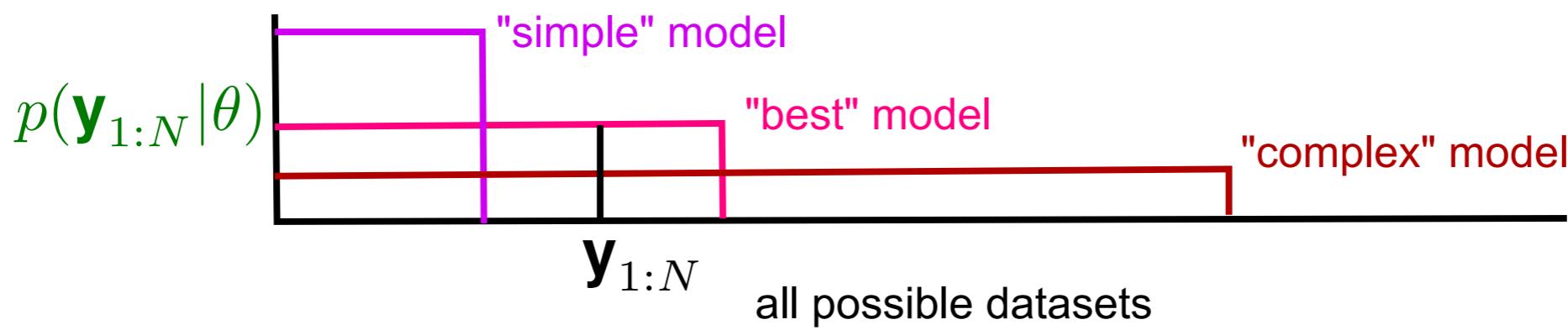
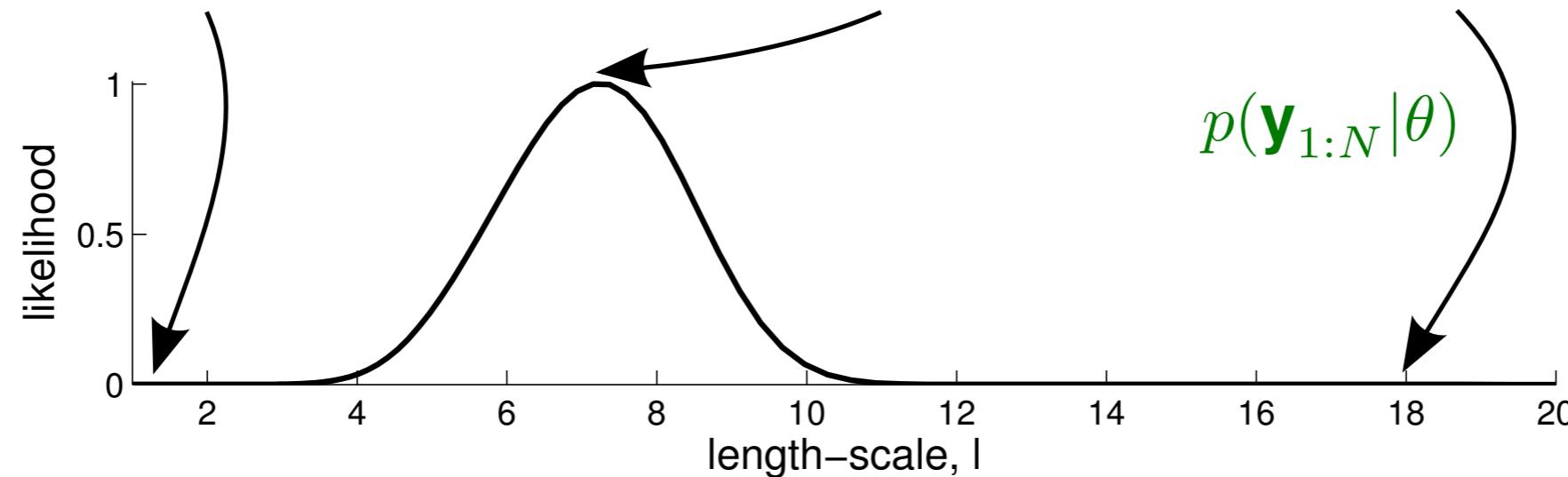
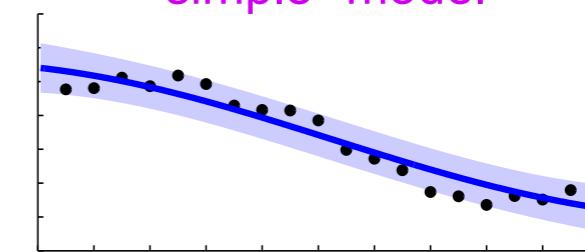
fits every training point  
"complex" model



"best" model



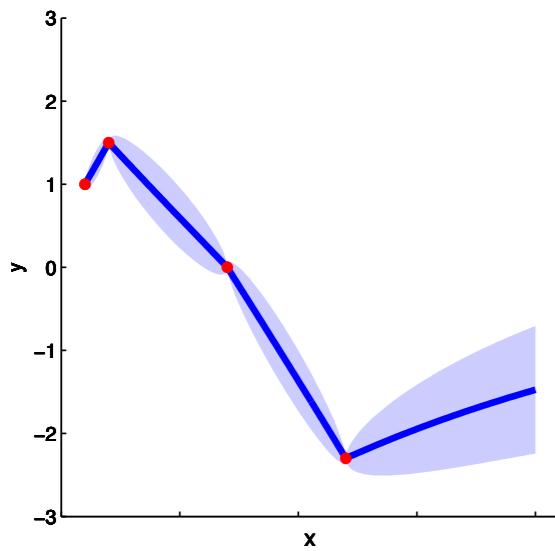
straight line-like  
"simple" model



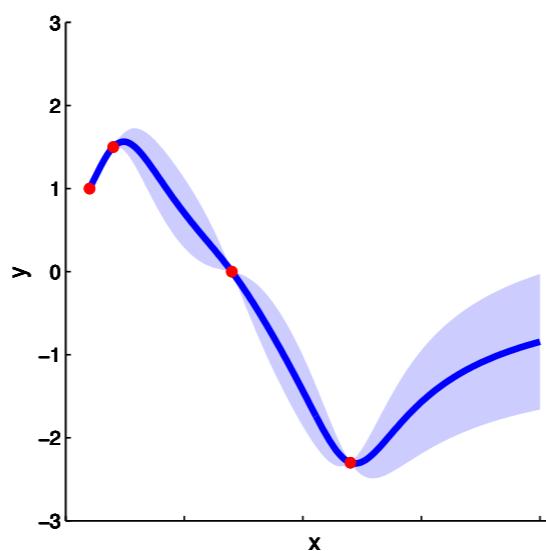
## Quick recap

### 5. GP kernels

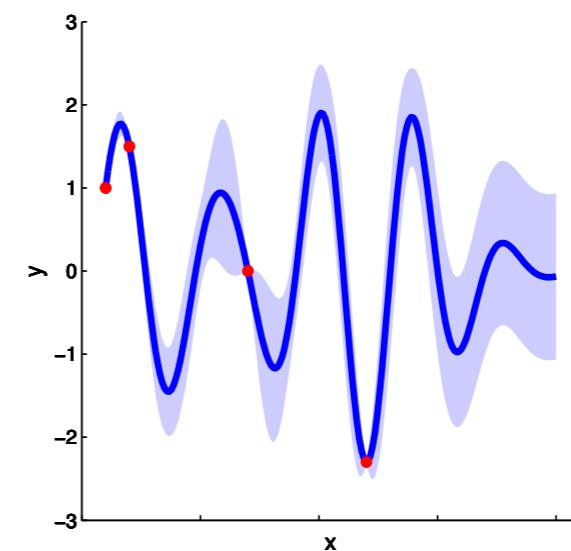
OU



RQ



periodic

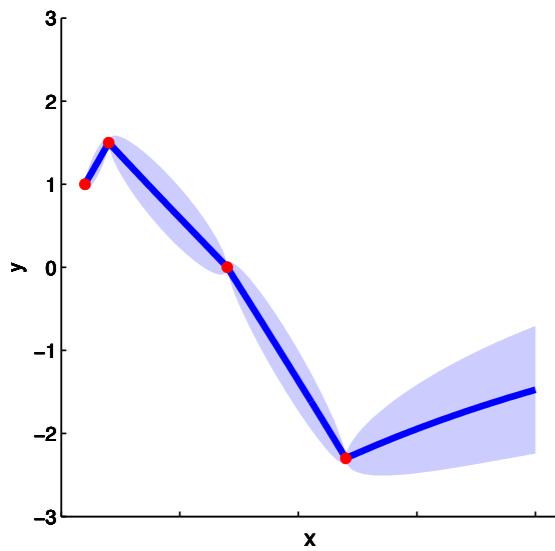


Choice of  
covariance function  
influence outcome

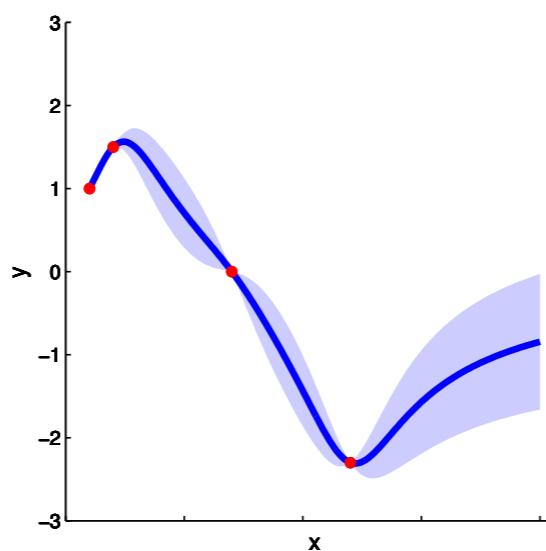
## Quick recap

### 5. GP kernels

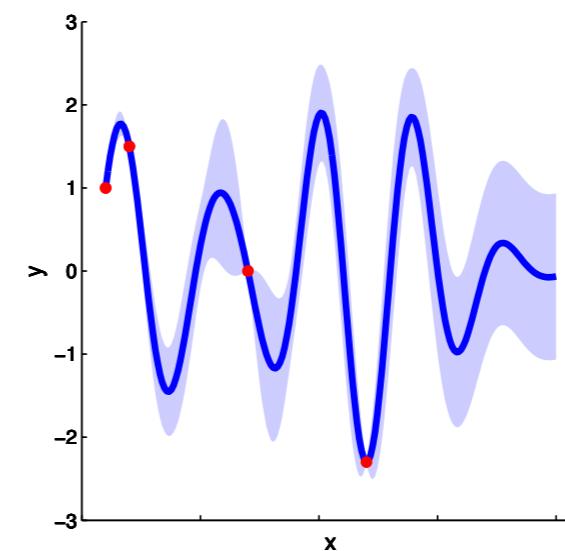
OU



RQ



periodic



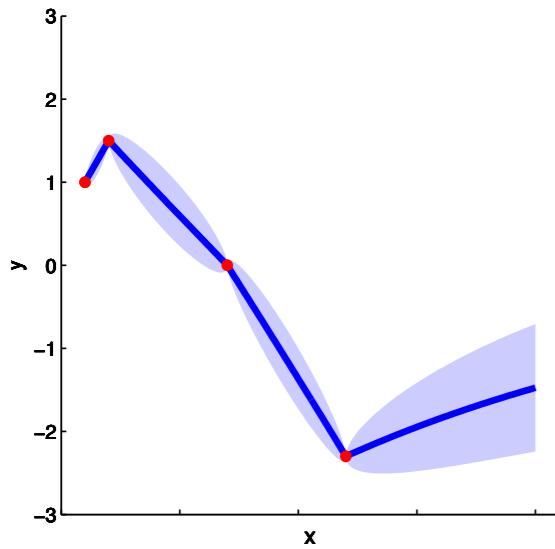
Choice of  
covariance function  
influence outcome

$$p(M|\mathbf{y}_{1:N}) = \frac{p(\mathbf{y}_{1:N}|M)p(M)}{\sum_{M'} p(\mathbf{y}_{1:N}|M')p(M')}$$

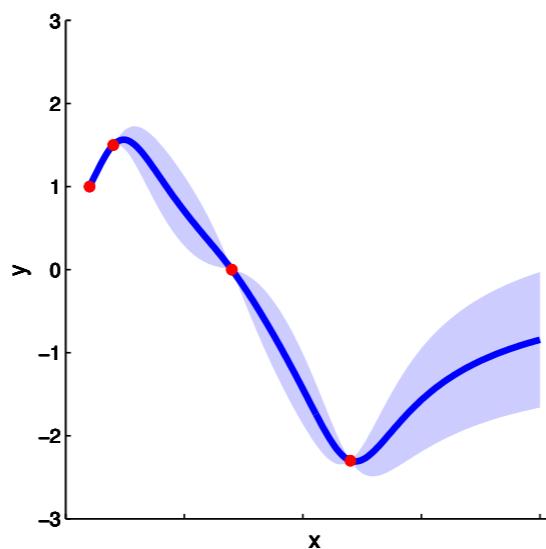
## Quick recap

### 5. GP kernels

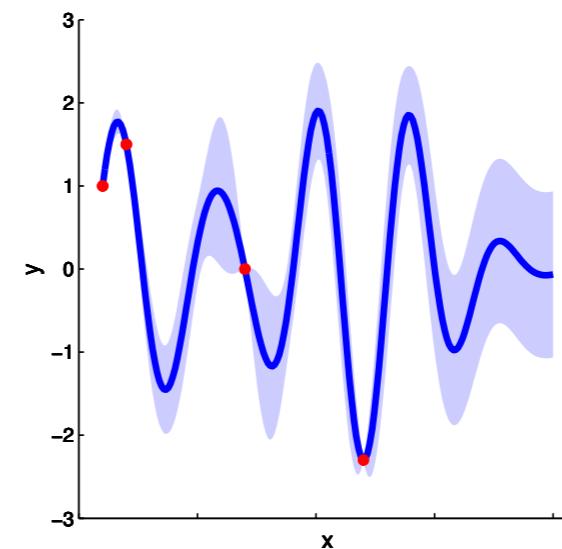
OU



RQ



periodic



Choice of  
covariance function  
influence outcome

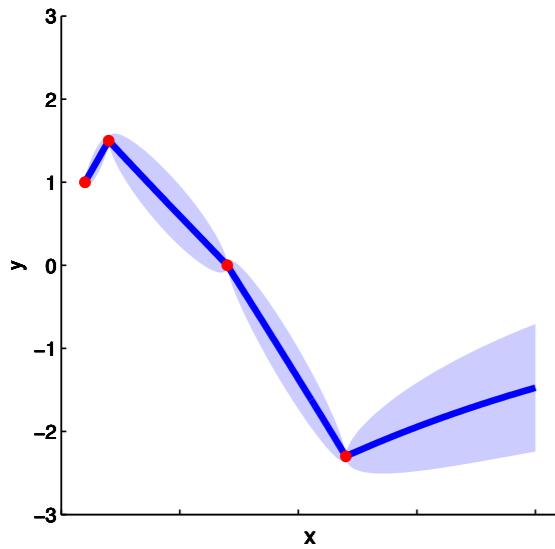
$$p(M|\mathbf{y}_{1:N}) = \frac{p(\mathbf{y}_{1:N}|M)p(M)}{\sum_{M'} p(\mathbf{y}_{1:N}|M')p(M')}$$

Prior over models

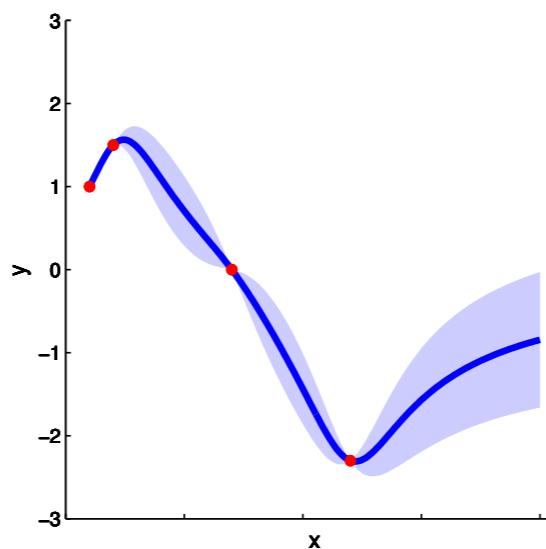
## Quick recap

### 5. GP kernels

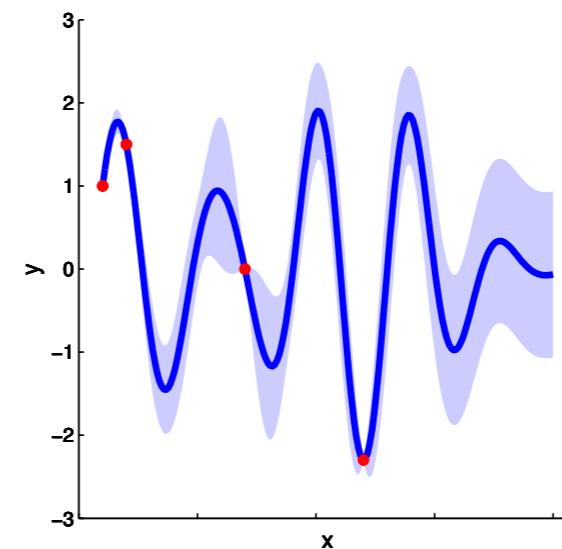
OU



RQ



periodic



Choice of  
covariance function  
influence outcome

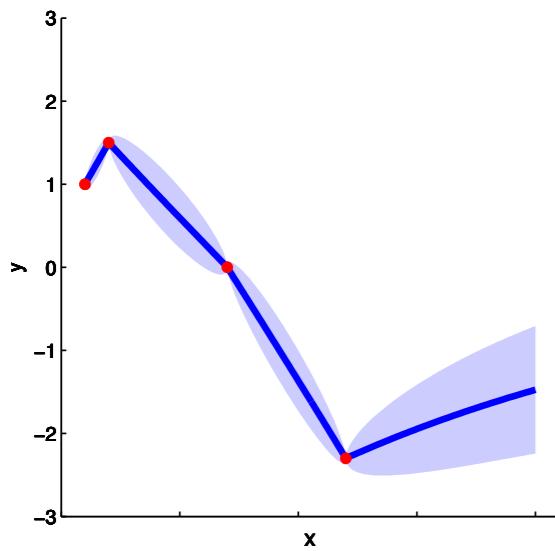
Marginal likelihood   Prior over models

$$p(M|\mathbf{y}_{1:N}) = \frac{p(\mathbf{y}_{1:N}|M)p(M)}{\sum_{M'} p(\mathbf{y}_{1:N}|M')p(M')}$$

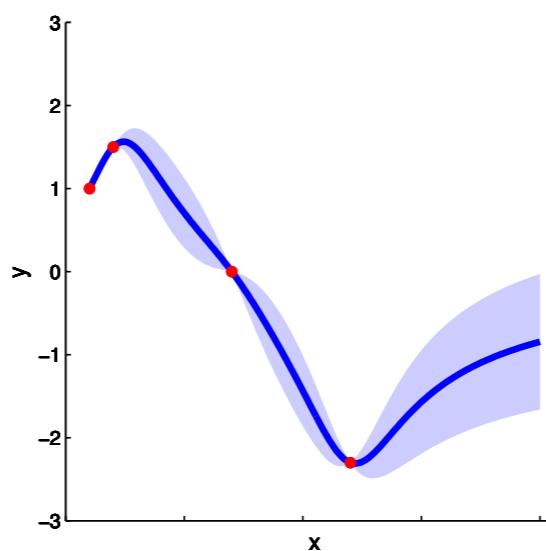
## Quick recap

### 5. GP kernels

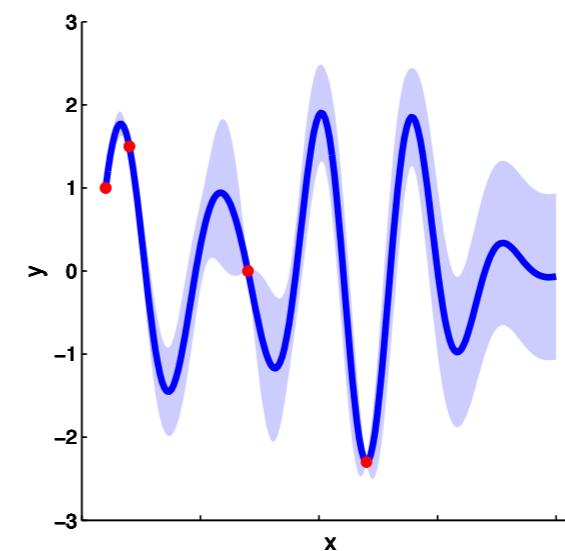
OU



RQ



periodic



Choice of covariance function influence outcome

Marginal likelihood   Prior over models

$$p(M|\mathbf{y}_{1:N}) = \frac{p(\mathbf{y}_{1:N}|M)p(M)}{\sum_{M'} p(\mathbf{y}_{1:N}|M')p(M')}$$

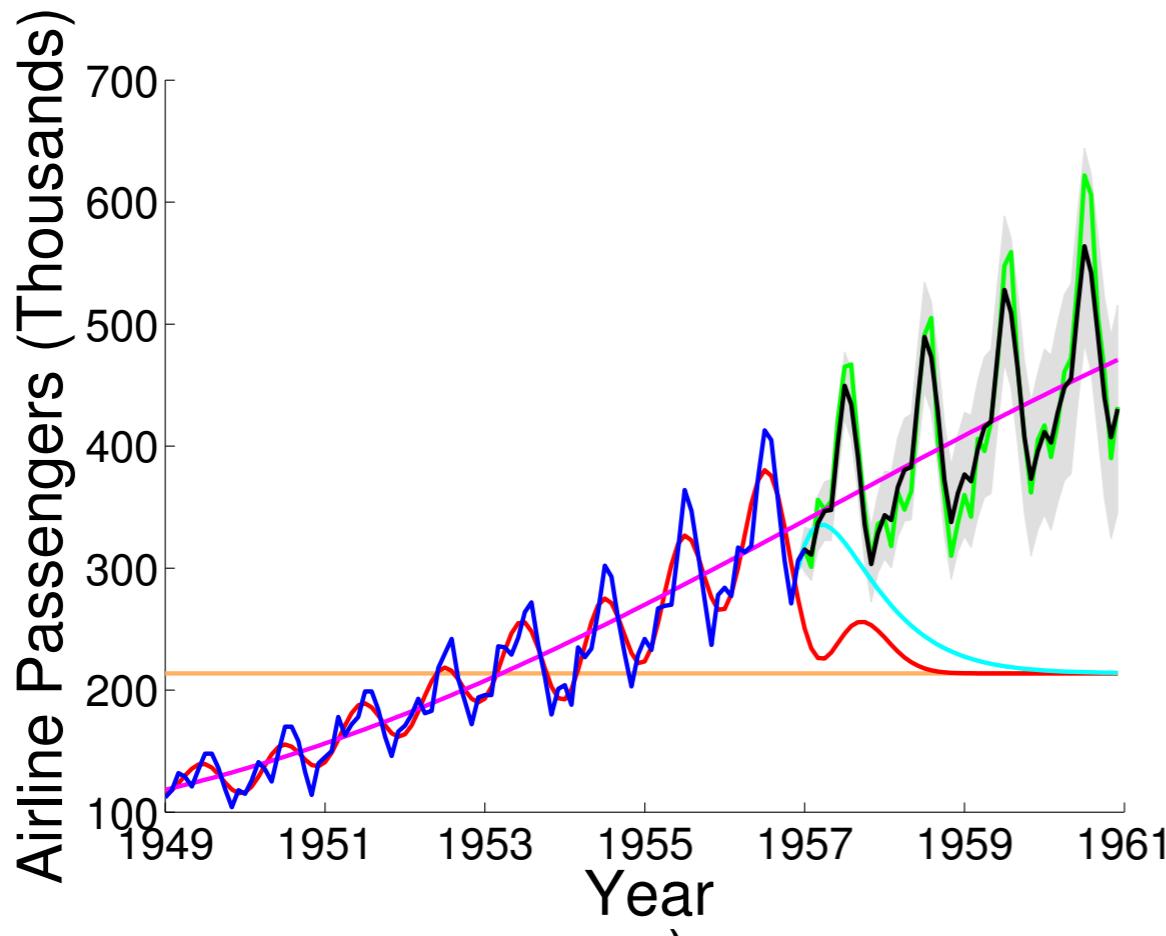
$$p(\mathbf{y}_{1:N}|M) = \int d\theta \, p(\mathbf{y}_{1:N}|\theta, M)p(\theta|M) \quad \text{Usually unpleasant}$$

## **Meta-kernel, learn which is right directly**

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$

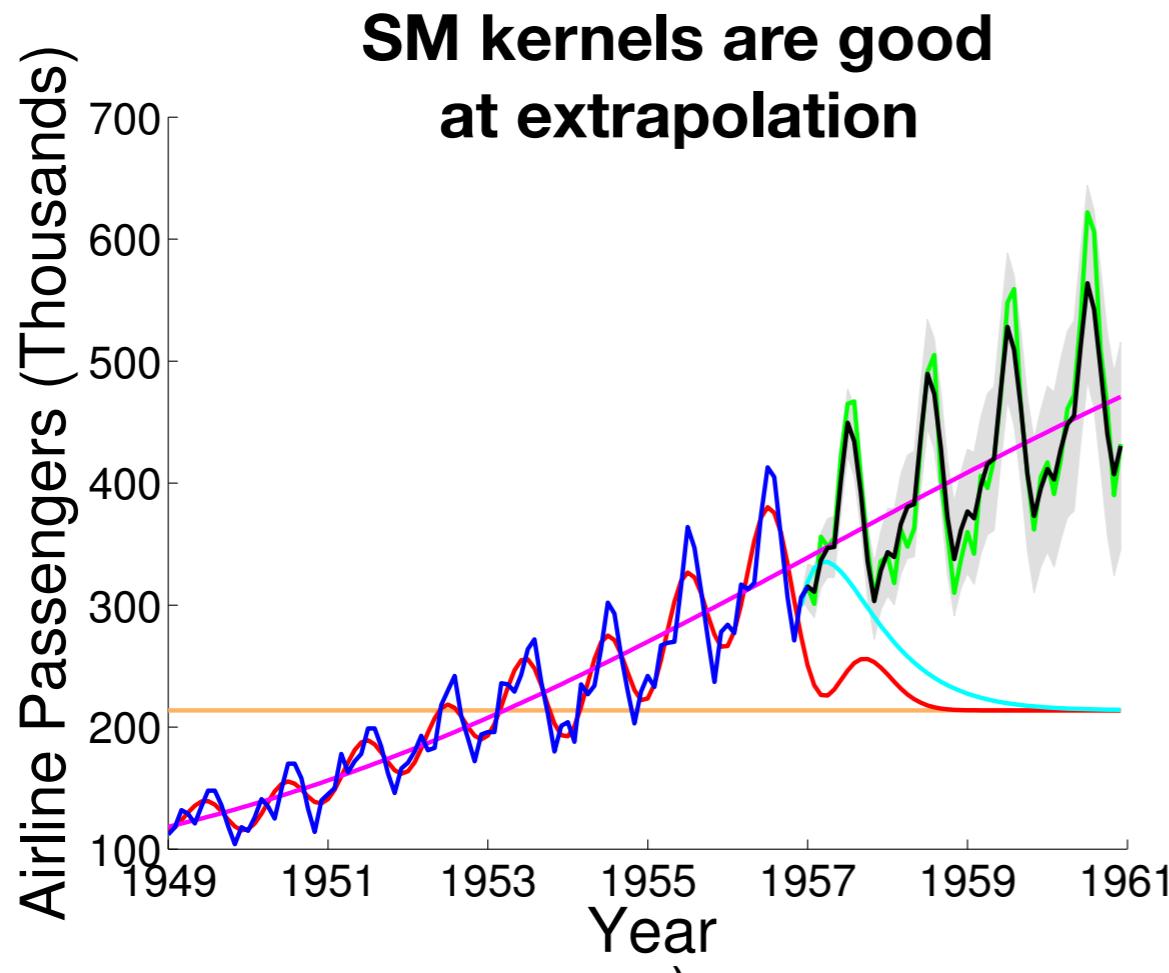
## Meta-kernel, learn which is right directly

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$



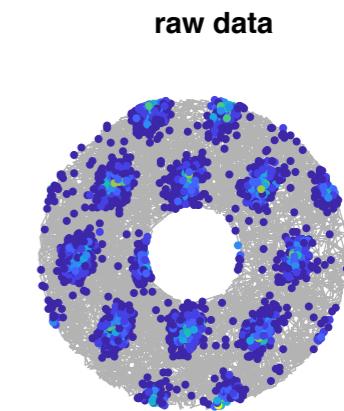
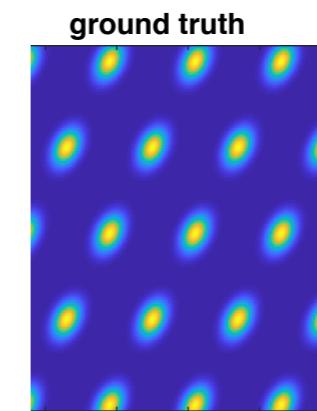
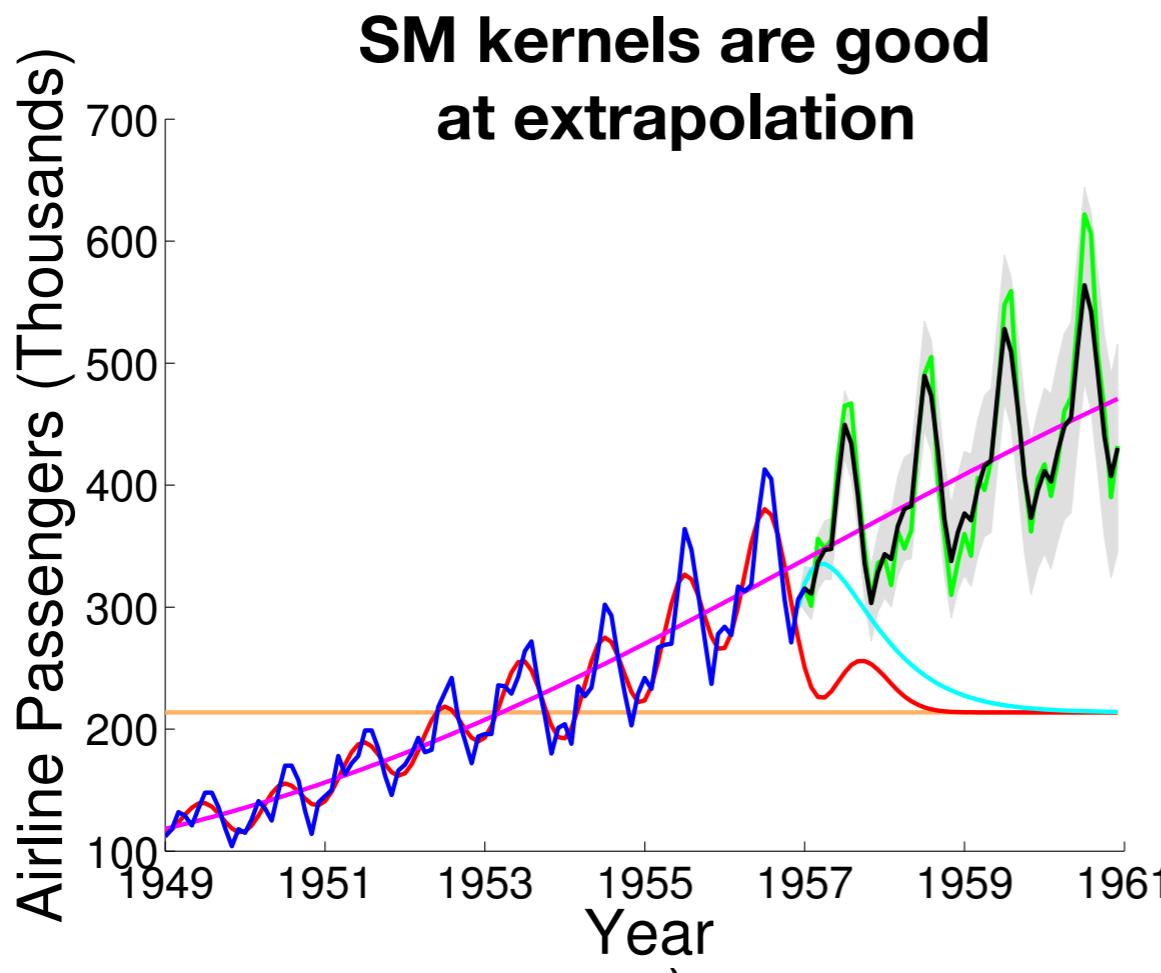
## Meta-kernel, learn which is right directly

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$



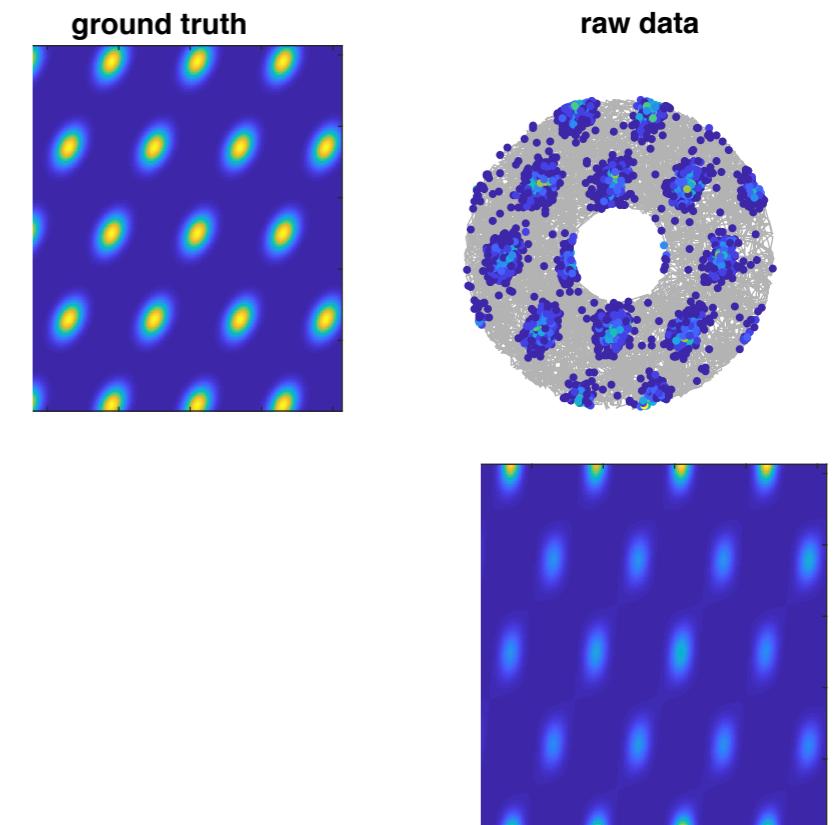
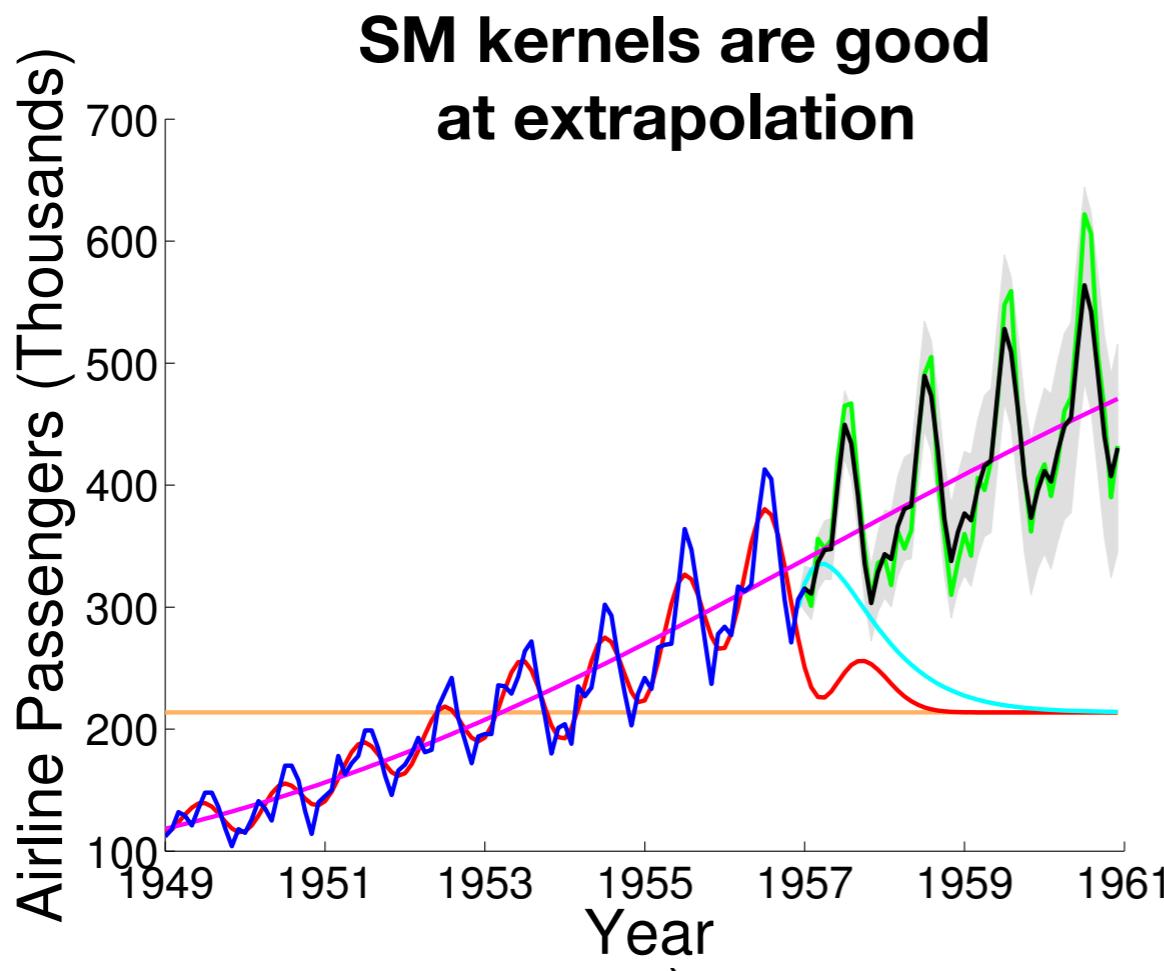
## Meta-kernel, learn which is right directly

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$



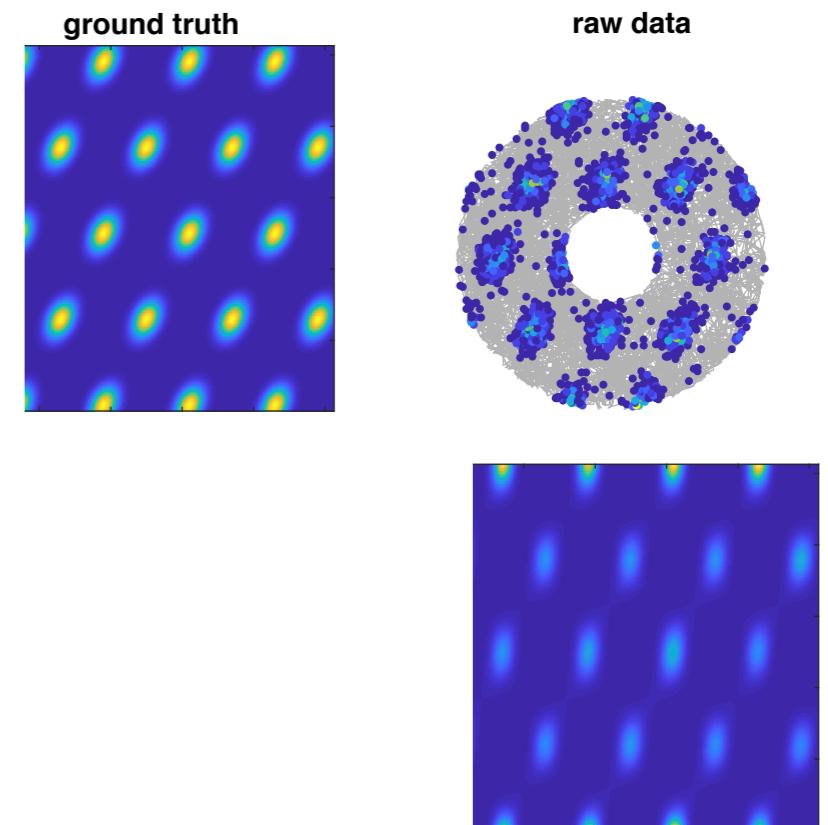
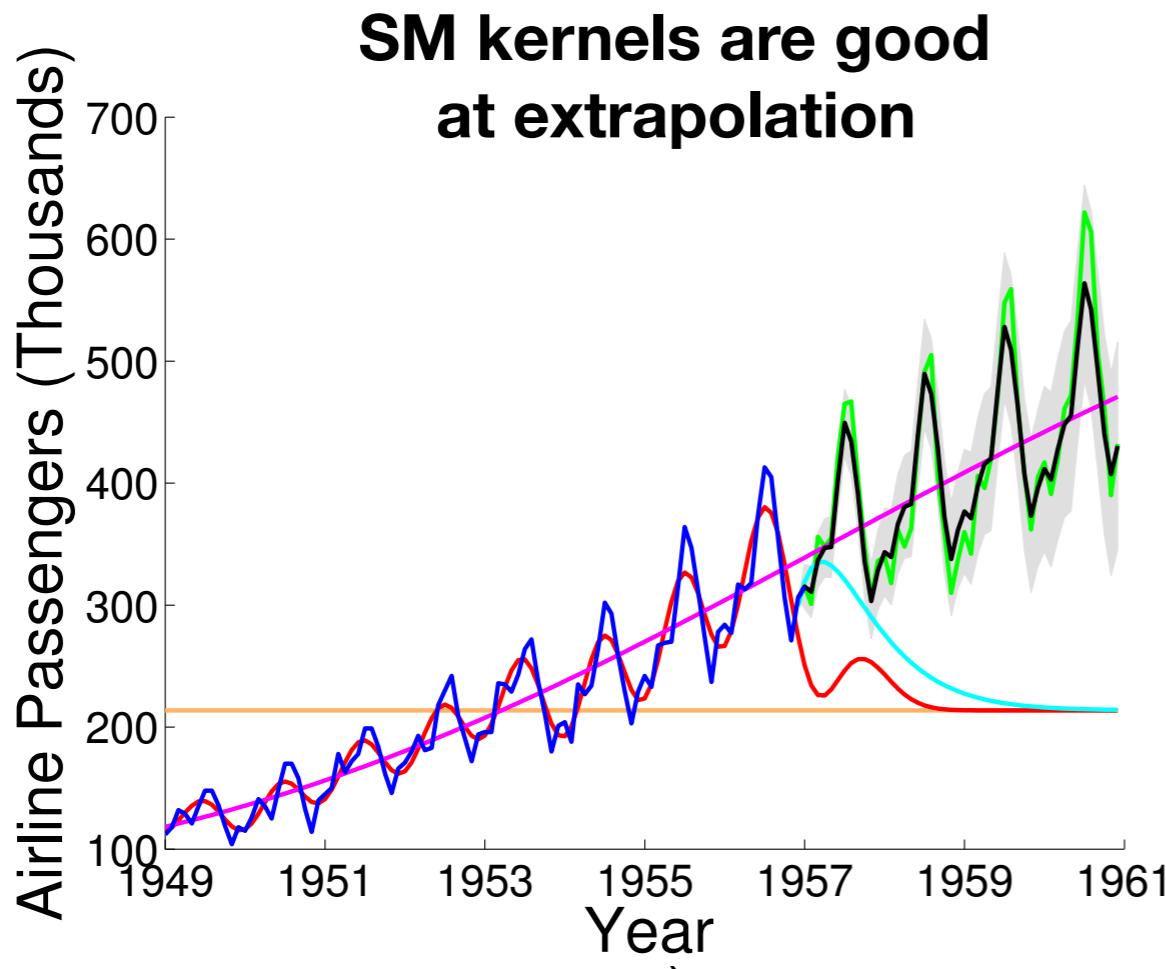
## Meta-kernel, learn which is right directly

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$



## Meta-kernel, learn which is right directly

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$

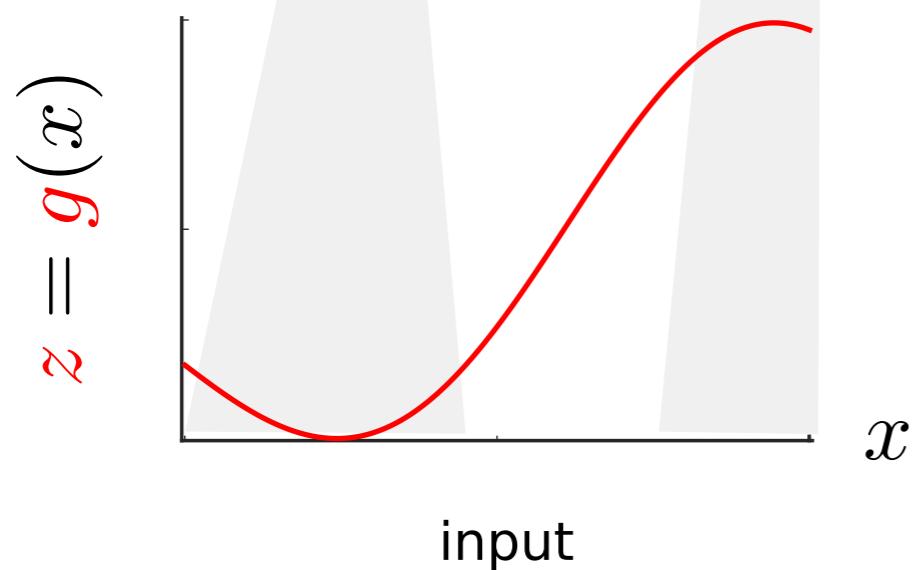
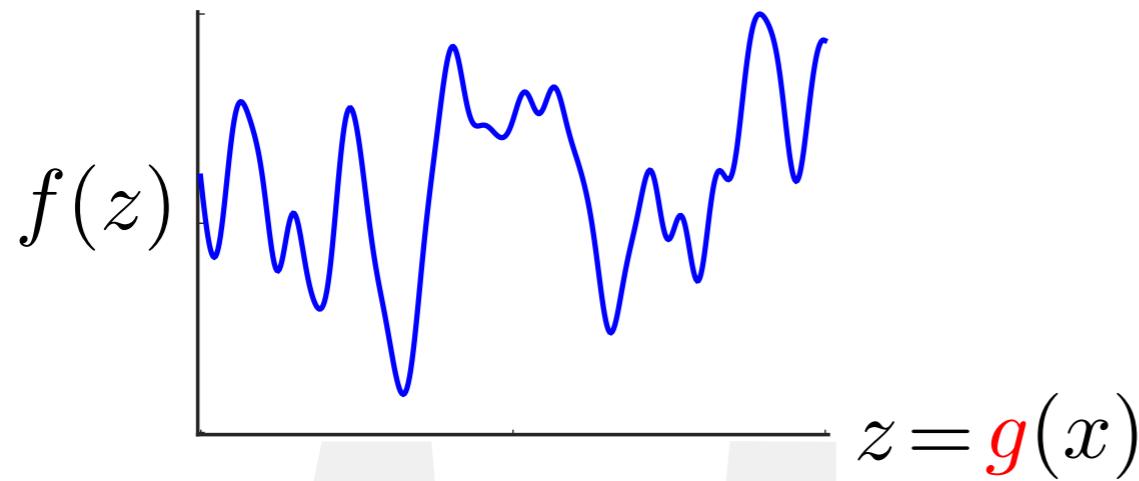


## Deep GPs

$$y(x) = \textcolor{blue}{f}(\textcolor{red}{g}(x)) + \sigma_y \epsilon$$

$$f(x) = \mathcal{GP}(0, K_f(x, x'))$$

$$g(x) = \mathcal{GP}(0, K_g(x, x'))$$

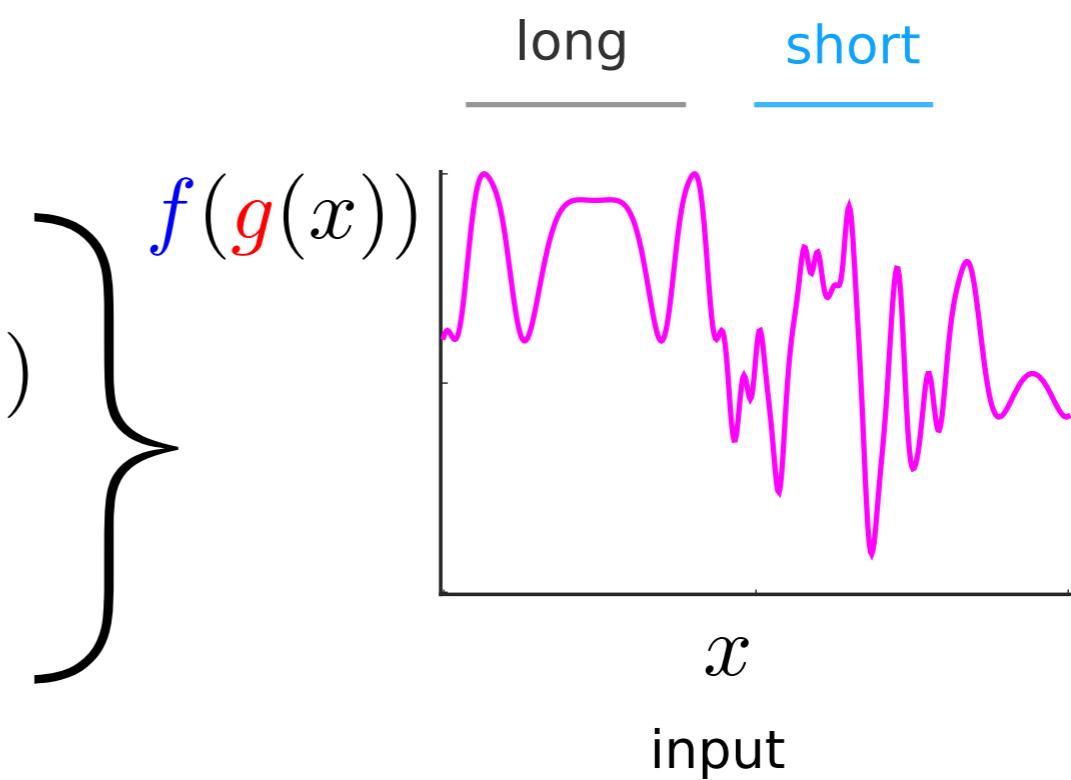
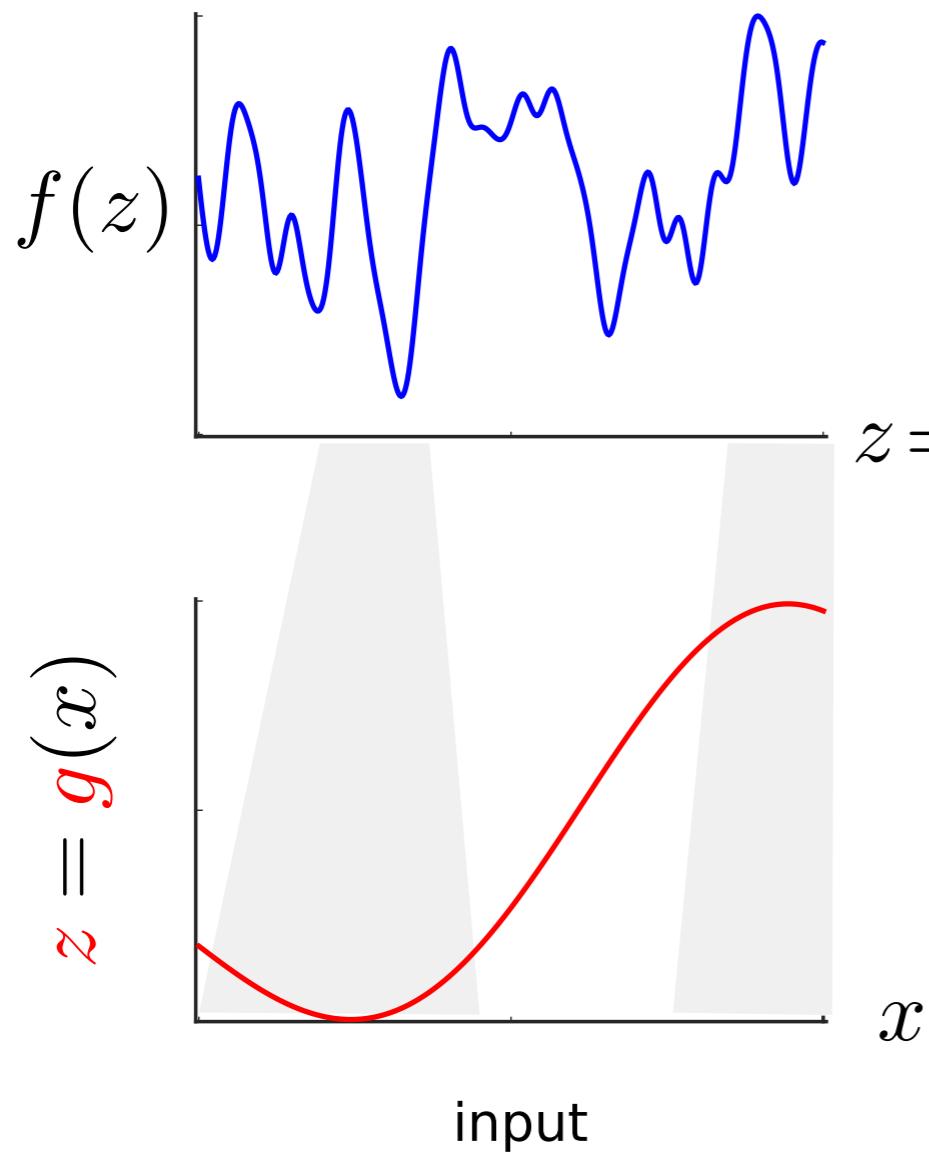


## Deep GPs

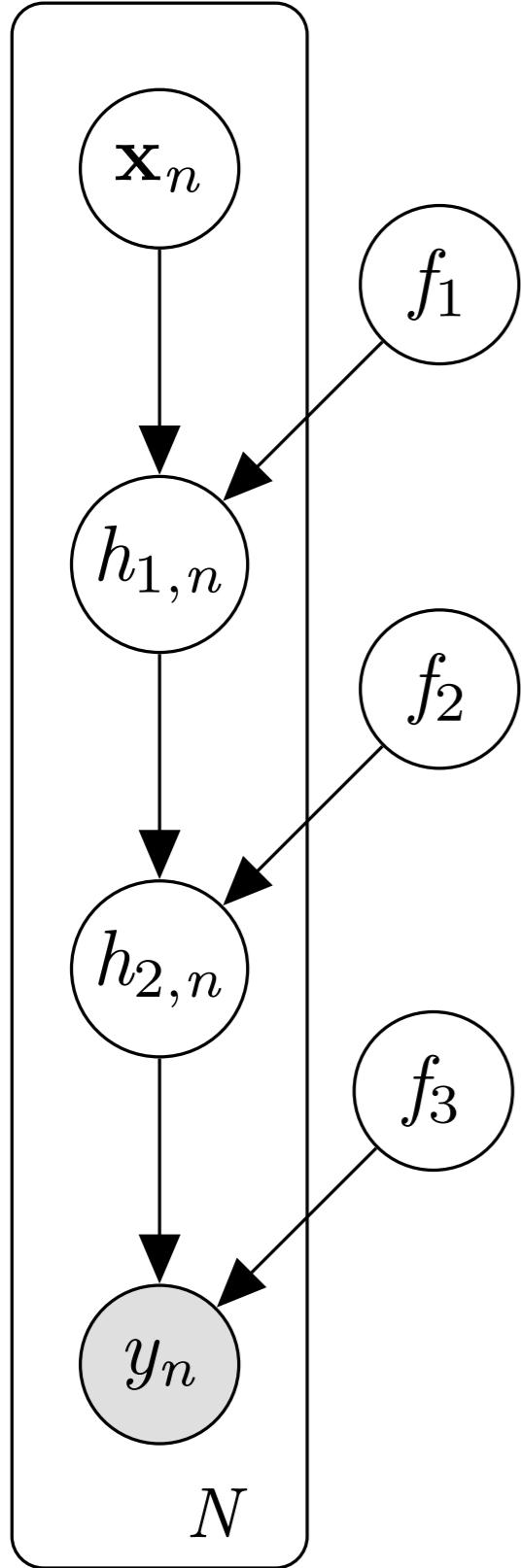
$$y(x) = \mathbf{f}(\mathbf{g}(x)) + \sigma_y \epsilon$$

$$\mathbf{f}(x) = \mathcal{GP}(0, K_f(x, x'))$$

$$\mathbf{g}(x) = \mathcal{GP}(0, K_g(x, x'))$$



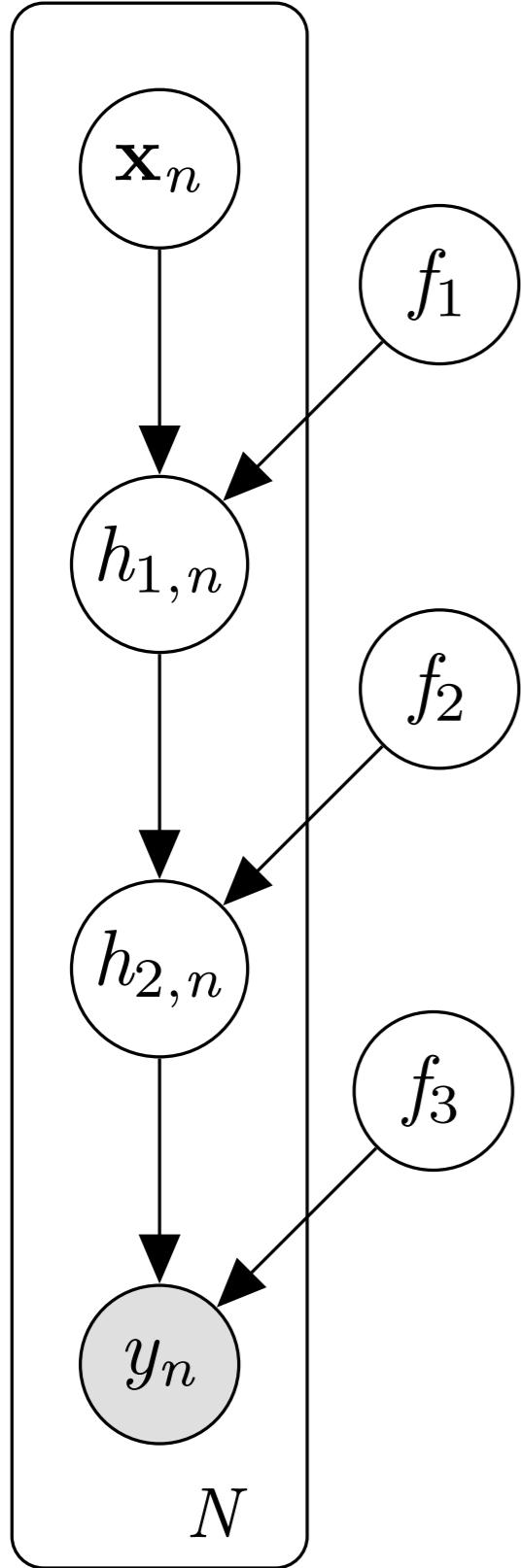
DGP = multi-layer neural network,  
infinitely wide hidden layer



$$f_l \sim \mathcal{GP}(0, k(., .))$$

$$y_n = g(\mathbf{x}_n) = f_L(f_{L-1}(\cdots f_2(f_1(\mathbf{x}_n)))) + \epsilon_n$$

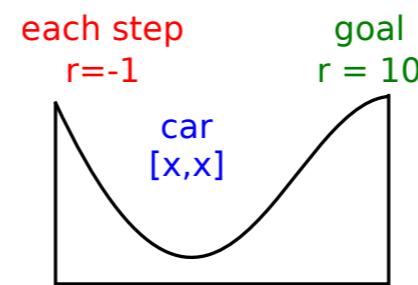
$$h_{L-1,n} := f_{L-1}(\cdots f_1(\mathbf{x}_n)), y_n = f_L(h_{L-1,n}) + \epsilon_n$$



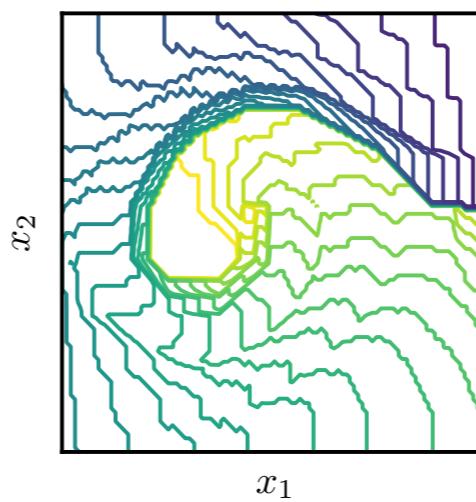
$$f_l \sim \mathcal{GP}(0, k(., .))$$

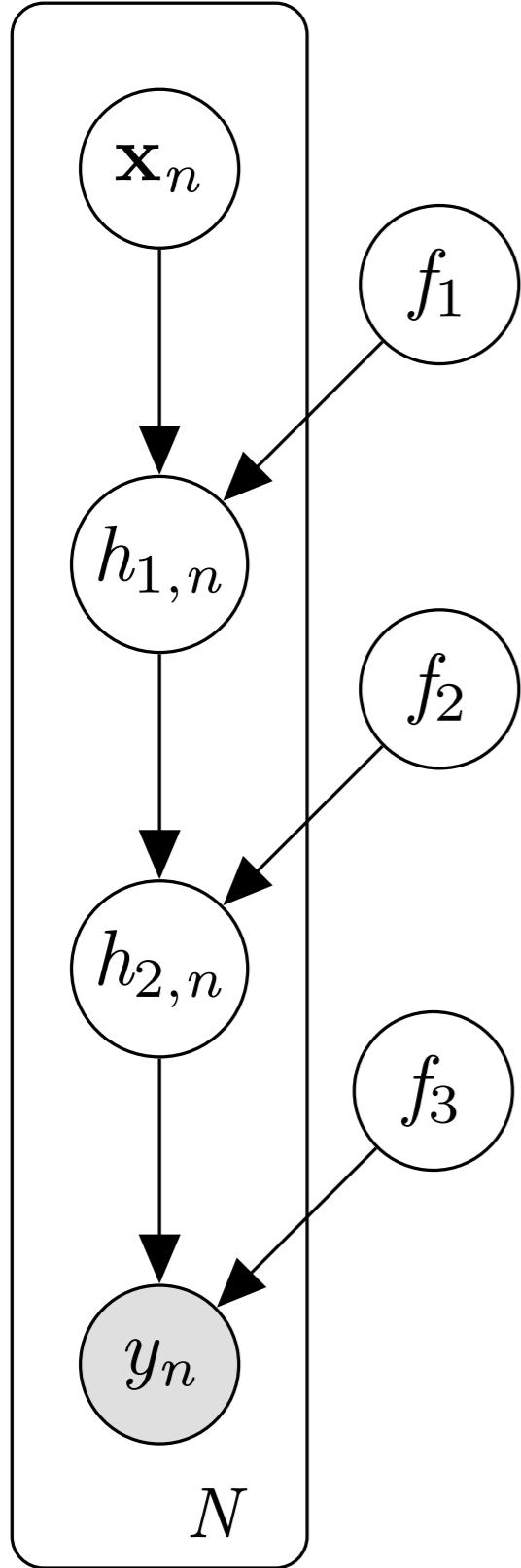
$$y_n = g(\mathbf{x}_n) = f_L(f_{L-1}(\cdots f_2(f_1(\mathbf{x}_n)))) + \epsilon_n$$

$$h_{L-1,n} := f_{L-1}(\cdots f_1(\mathbf{x}_n)), y_n = f_L(h_{L-1,n}) + \epsilon_n$$



Value function

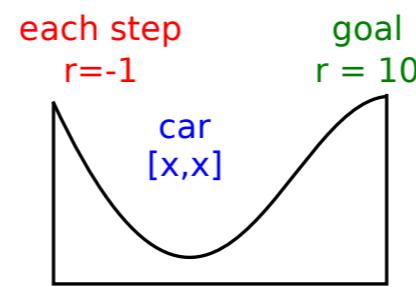




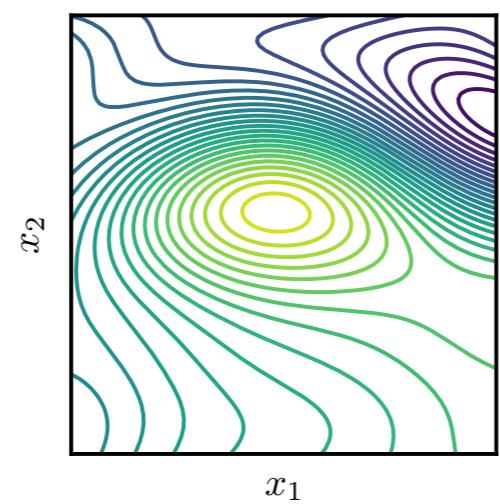
$$f_l \sim \mathcal{GP}(0, k(., .))$$

$$y_n = g(\mathbf{x}_n) = f_L(f_{L-1}(\cdots f_2(f_1(\mathbf{x}_n)))) + \epsilon_n$$

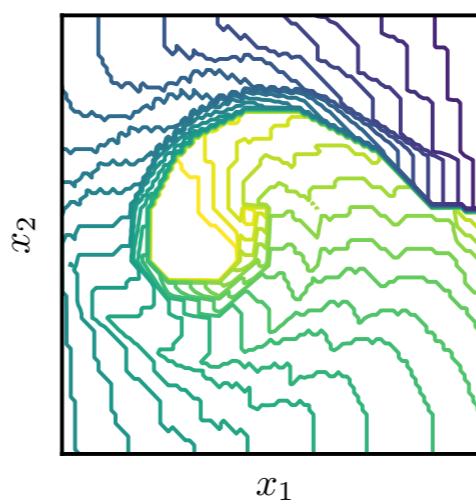
$$h_{L-1,n} := f_{L-1}(\cdots f_1(\mathbf{x}_n)), y_n = f_L(h_{L-1,n}) + \epsilon_n$$

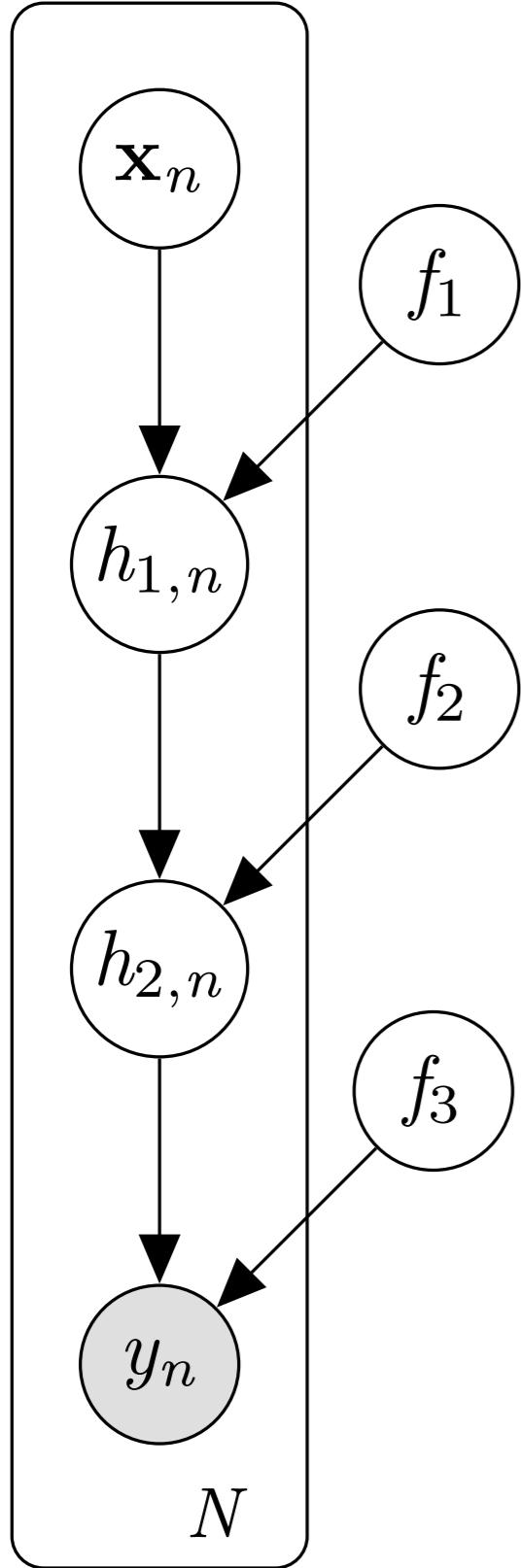


GP fit



Value function

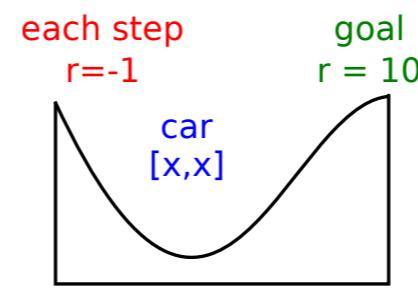




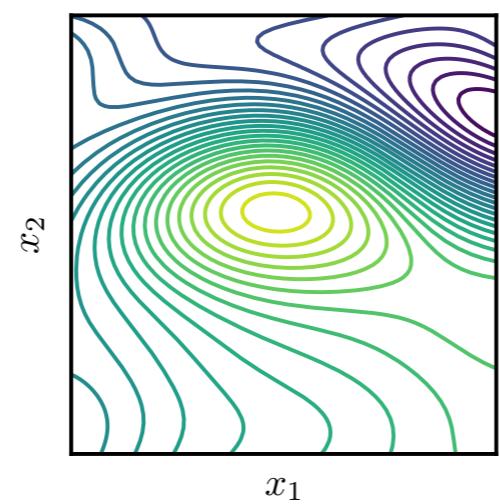
$$f_l \sim \mathcal{GP}(0, k(., .))$$

$$y_n = g(\mathbf{x}_n) = f_L(f_{L-1}(\cdots f_2(f_1(\mathbf{x}_n)))) + \epsilon_n$$

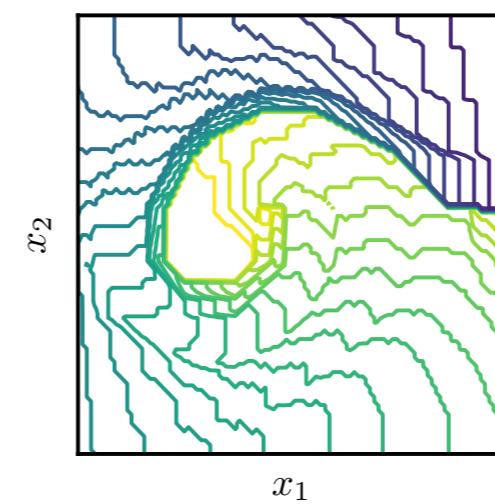
$$h_{L-1,n} := f_{L-1}(\cdots f_1(\mathbf{x}_n)), y_n = f_L(h_{L-1,n}) + \epsilon_n$$



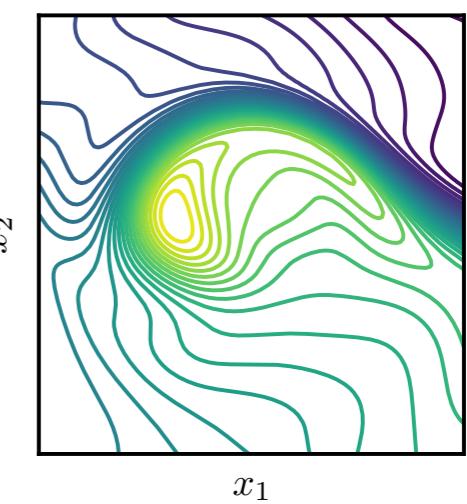
GP fit



Value function



DGP fit



## **What are Gaussian Processes good for?**

## **What are Gaussian Processes good for?**

### **Strengths**

## What are Gaussian Processes good for?

### Strengths

**interpretable:** covariance functions specify easy-to-explain high-level properties of functions

## What are Gaussian Processes good for?

### Strengths

**interpretable:** covariance functions specify easy-to-explain high-level properties of functions

**data-efficient:** non-parametric + Bayesian  $\Rightarrow$  lots of flexibility + avoid overfitting

## What are Gaussian Processes good for?

### Strengths

**interpretable:** covariance functions specify easy-to-explain high-level properties of functions

**data-efficient:** non-parametric + Bayesian  $\Rightarrow$  lots of flexibility + avoid overfitting

**optimal decision making:** well-calibrated uncertainties: knows when it does not know

## What are Gaussian Processes good for?

### Strengths

**interpretable:** covariance functions specify easy-to-explain high-level properties of functions

**data-efficient:** non-parametric + Bayesian  $\Rightarrow$  lots of flexibility + avoid overfitting

**optimal decision making:** well-calibrated uncertainties: knows when it does not know

### Weaknesses

## What are Gaussian Processes good for?

### Strengths

**interpretable:** covariance functions specify easy-to-explain high-level properties of functions

**data-efficient:** non-parametric + Bayesian  $\Rightarrow$  lots of flexibility + avoid overfitting

**optimal decision making:** well-calibrated uncertainties: knows when it does not know

### Weaknesses

**large datasets:**  $N \leq 10^5$  unless there is special structure (invert and store covariance matrix)

## What are Gaussian Processes good for?

### Strengths

**interpretable:** covariance functions specify easy-to-explain high-level properties of functions

**data-efficient:** non-parametric + Bayesian  $\Rightarrow$  lots of flexibility + avoid overfitting

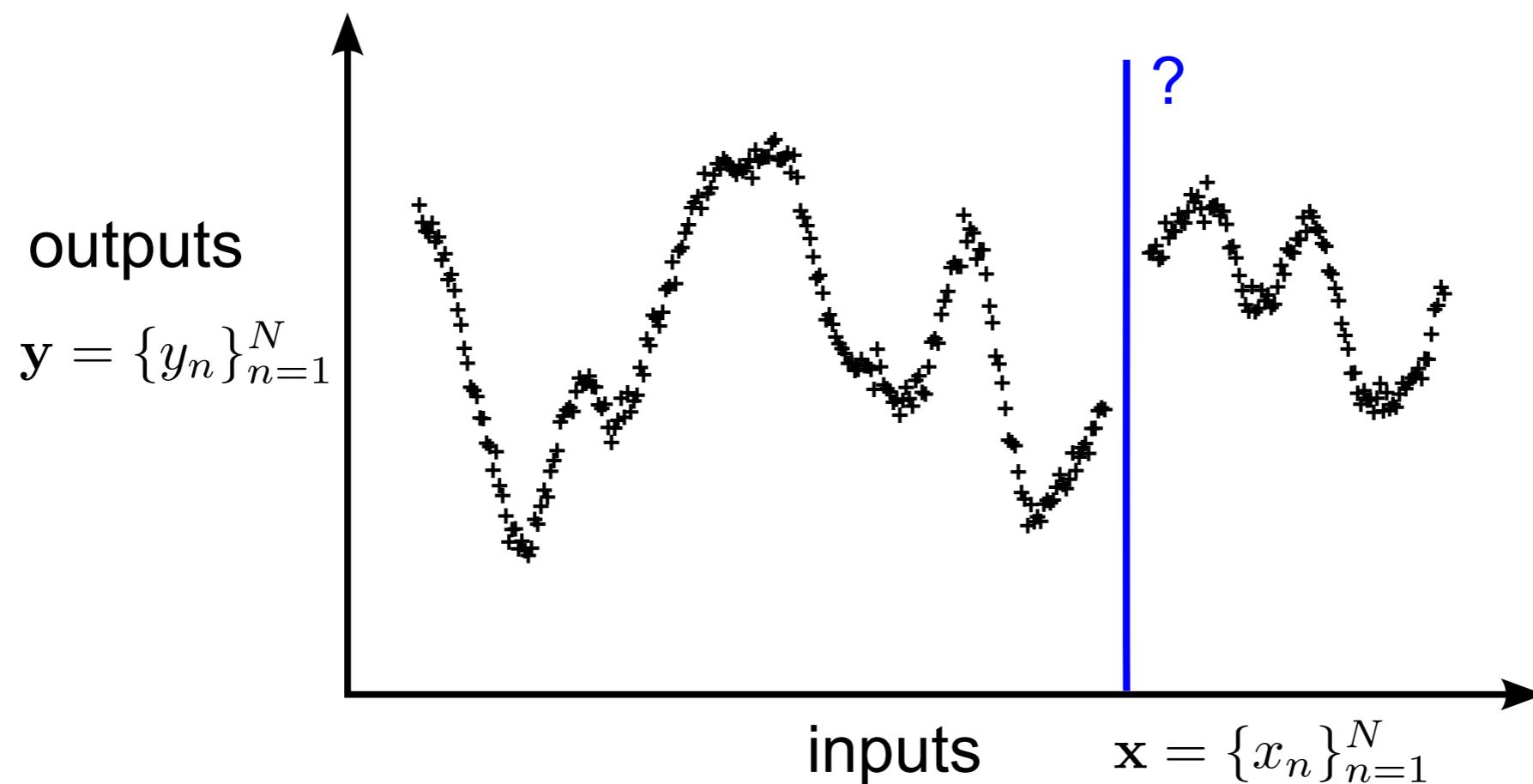
**optimal decision making:** well-calibrated uncertainties: knows when it does not know

### Weaknesses

**large datasets:**  $N \leq 10^5$  unless there is special structure (invert and store covariance matrix)

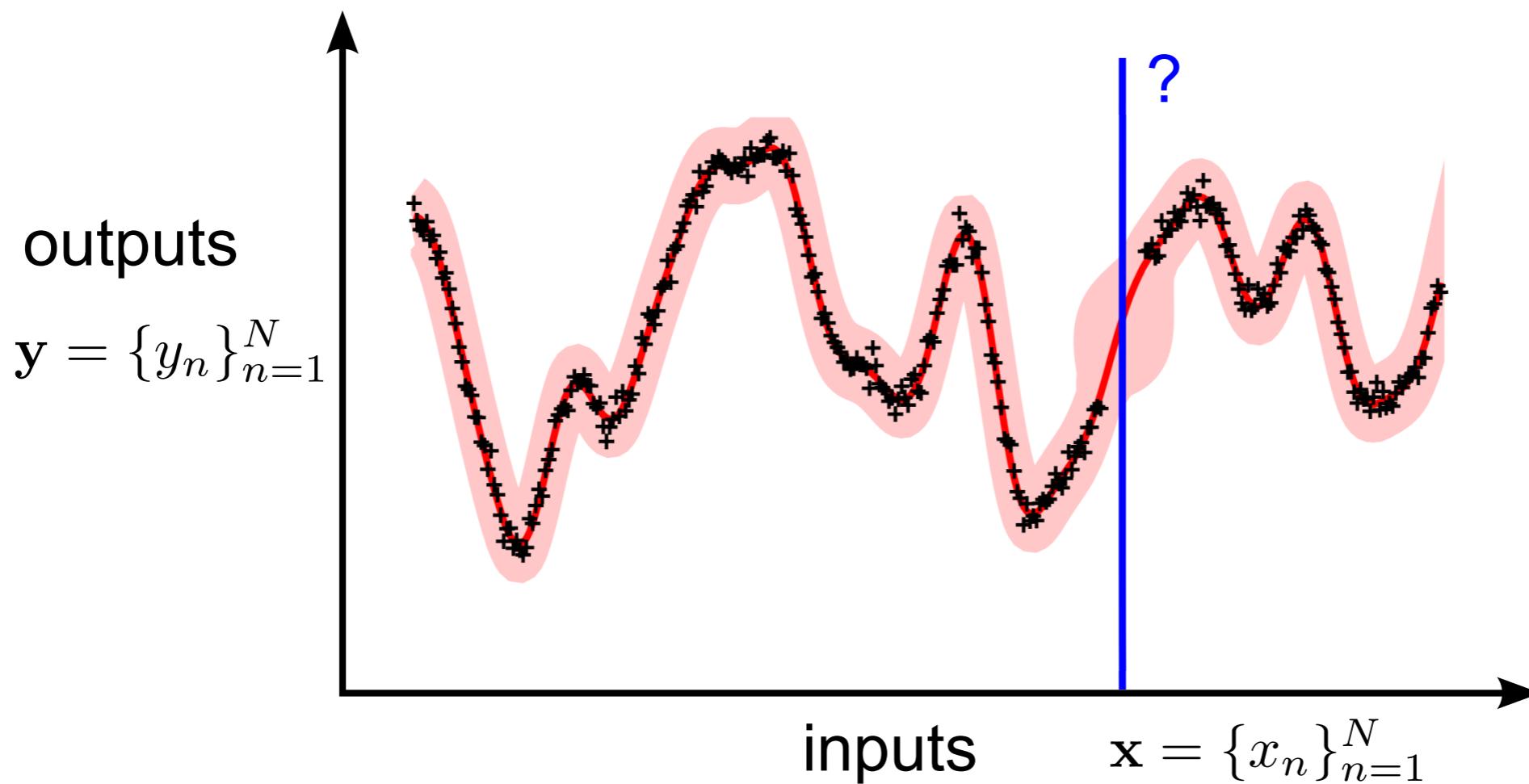
**high-dimensional inputs spaces:**  $D \leq 10^2$  unless there is special structure, e.g. Kronecker (compute pair-wise elements of covariance function)

## How do me make these work in practice?



## How do me make these work in practice?

$$\begin{array}{ccc} p(f|\theta) = \mathcal{GP}(f; 0, K_\theta) & \xrightarrow{\text{inference \& learning}} & p(f|\mathbf{y}, \mathbf{x}, \theta) \\ p(y_n|f, x_n, \theta) & & p(\mathbf{y}|\mathbf{x}, \theta) \end{array}$$

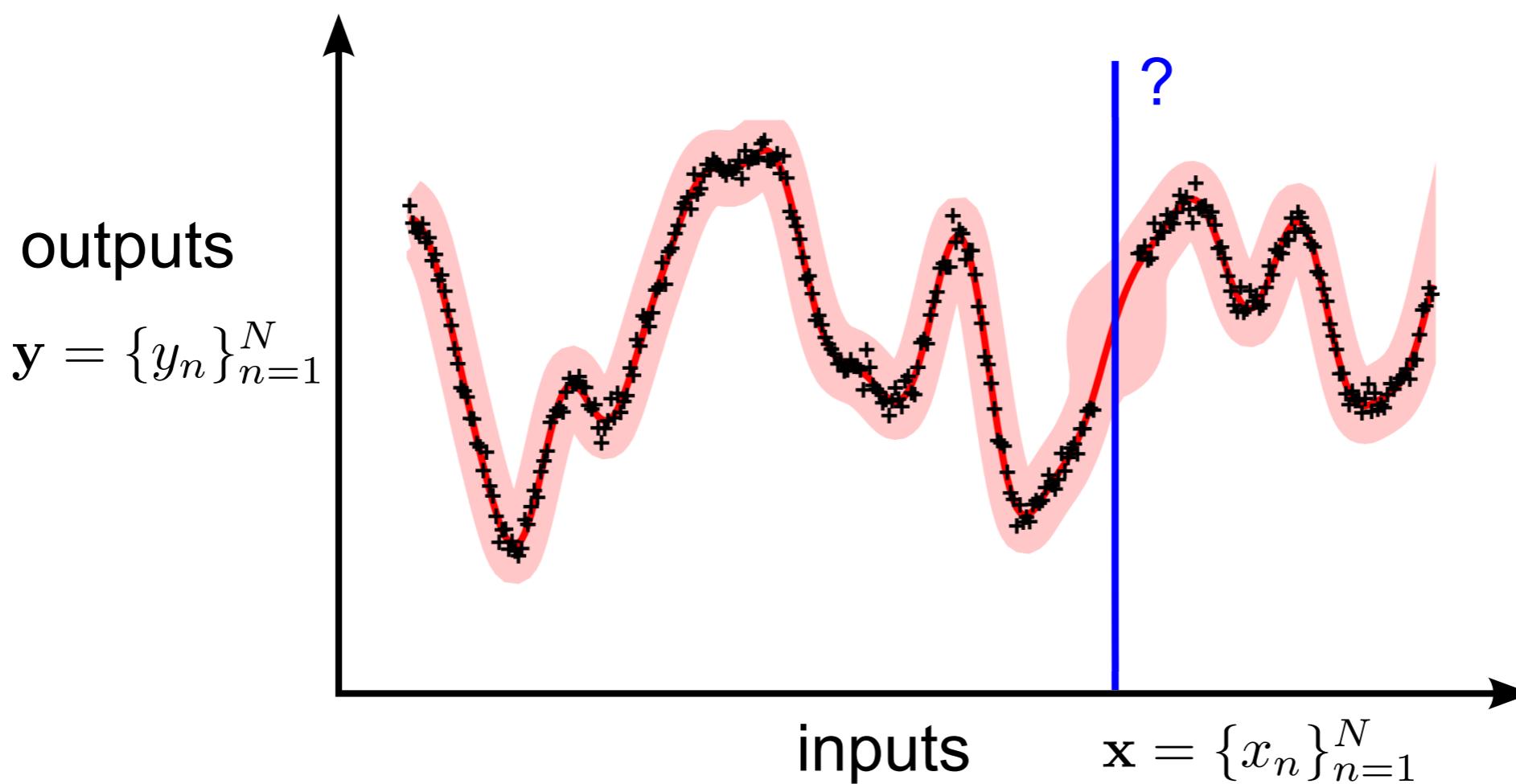


## How do me make these work in practice?

$$p(f|\theta) = \mathcal{GP}(f; 0, K_\theta) \quad \xrightarrow{\text{inference & learning}} \quad p(f|\mathbf{y}, \mathbf{x}, \theta)$$

$$p(y_n|f, x_n, \theta) \quad \xrightarrow{\text{intractabilities}} \quad p(\mathbf{y}|\mathbf{x}, \theta)$$

computational  $\mathcal{O}(N^3)$   
analytic



$$p(f|\theta) = \mathcal{GP}(f; 0, K_\theta)$$

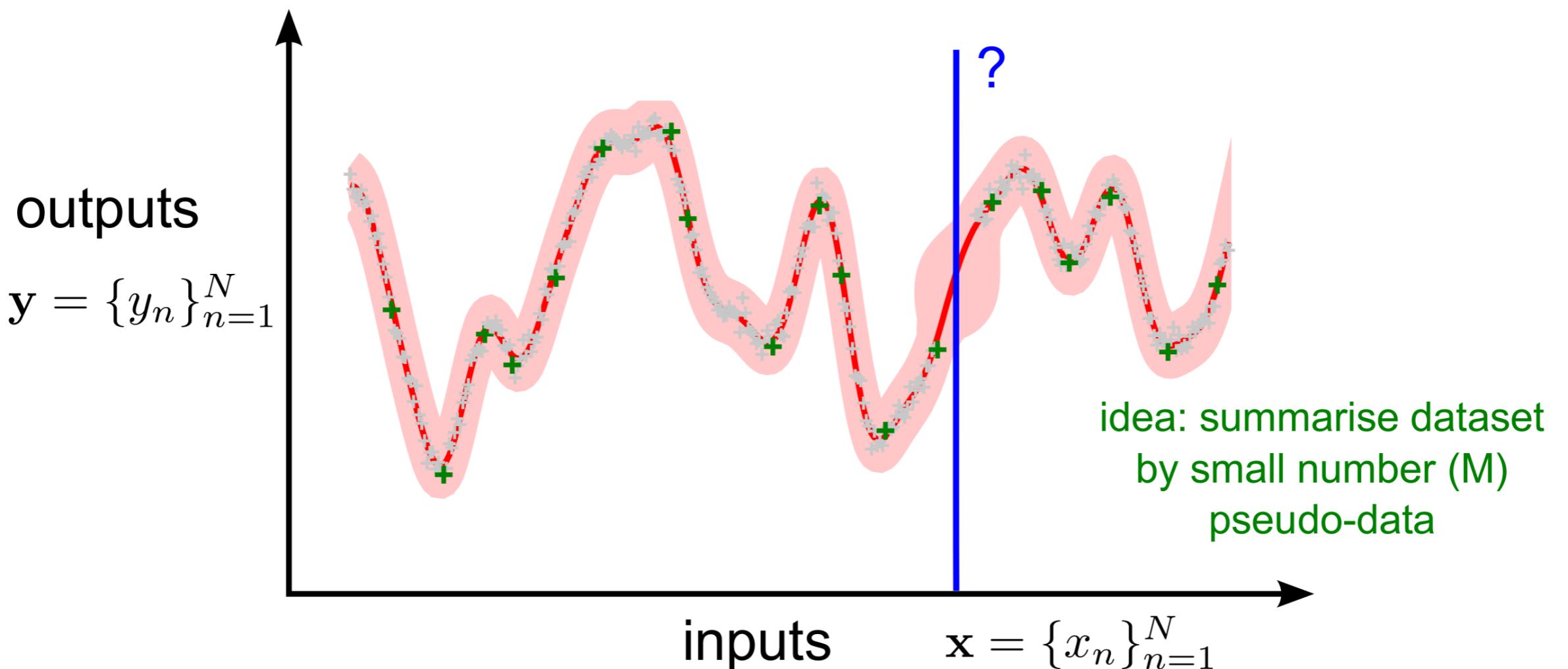
$$p(y_n|f, x_n, \theta)$$

inference & learning

intractabilities  
computational  $\mathcal{O}(N^3)$   
analytic

$$p(f|\mathbf{y}, \mathbf{x}, \theta)$$

$$p(\mathbf{y}|\mathbf{x}, \theta)$$



# A Brief History of Gaussian Process Approximations

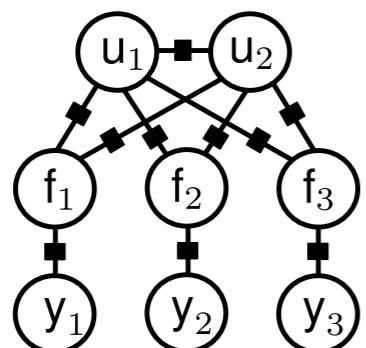
approximate generative model  
exact inference

methods employing  
pseudo-data

exact generative model  
approximate inference

$$\text{div}[p(\mathbf{f}, \mathbf{y}) || q(\mathbf{f}, \mathbf{y})]$$

A Unifying View of Sparse  
Approximate Gaussian  
Process Regression  
Quinonero-Candela &  
Rasmussen, 2005  
(FITC, PITC, DTC)



$$\text{div}[p(\mathbf{f}|\mathbf{y}) || q(\mathbf{f})]$$

FITC  
PITC  
DTC

VFE  
EP  
PP

A Unifying Framework for  
Sparse Gaussian Process  
Approximation using  
Power Expectation  
Propagation  
Bui, Yan and Turner, 2016  
(VFE, EP, FITC, PITC ...)

FITC: Snelson et al. "Sparse Gaussian Processes using Pseudo-inputs"

PITC: Snelson et al. "Local and global sparse Gaussian process approximations"

EP: Csato and Opper 2002 / Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

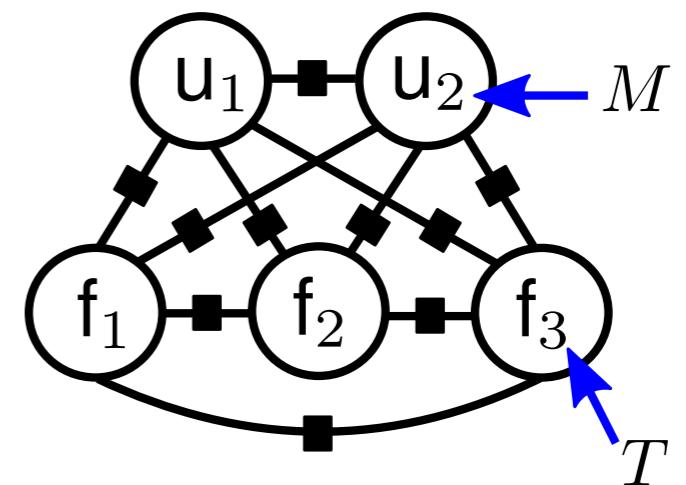
VFE: Titsias "Variational Learning of Inducing Variables in Sparse Gaussian Processes"

DTC / PP: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

# FITC

## 1. Add pseudo points

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_{ff} & K_{fu} \\ K_{uf} & K_{uu} \end{bmatrix} \right)$$

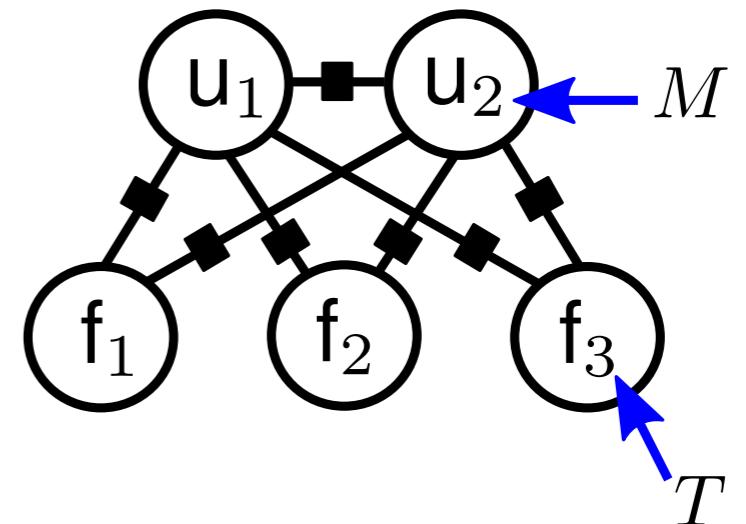


# FITC

## 1. Add pseudo points

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_{ff} & K_{fu} \\ K_{uf} & K_{uu} \end{bmatrix} \right)$$

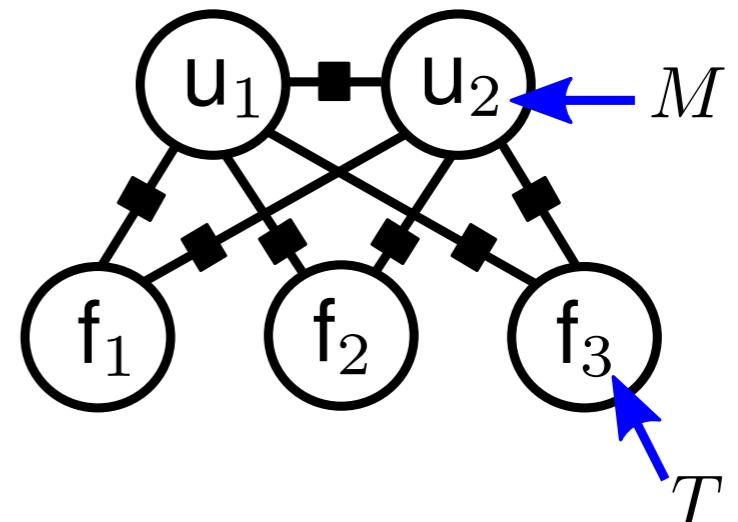
## 2. drop direct f-f dependencies



# FITC

## 1. Add pseudo points

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_{ff} & K_{fu} \\ K_{uf} & K_{uu} \end{bmatrix} \right)$$



## 2. drop direct f-f dependencies

## 3. calibrate model

(e.g. using KL divergence, many choices)

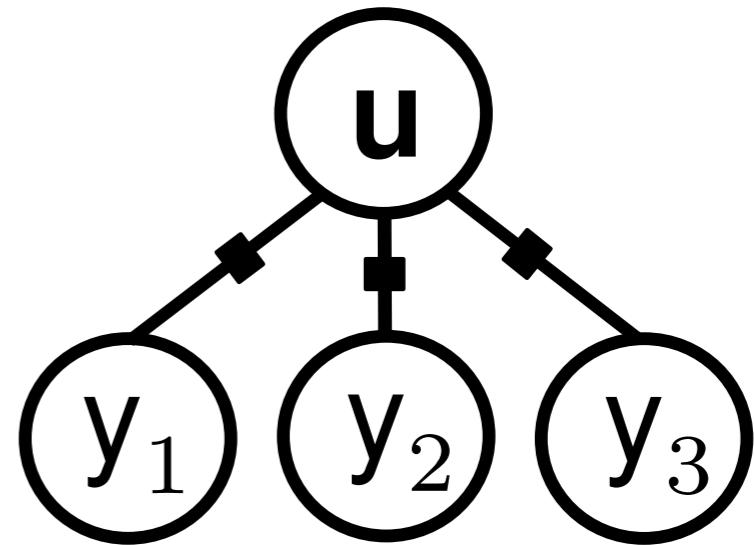
$$\arg \min_{q(\mathbf{u}), \{q(\mathbf{f}_t | \mathbf{u})\}_{t=1}^T} \text{KL}(p(\mathbf{f}, \mathbf{u}) || q(\mathbf{u}) \prod_{t=1}^T q(\mathbf{f}_t | \mathbf{u})) \implies \begin{aligned} q(\mathbf{u}) &= p(\mathbf{u}) \\ q(\mathbf{f}_t | \mathbf{u}) &= p(\mathbf{f}_t | \mathbf{u}) \end{aligned}$$

equal to exact conditionals

# FITC

## 1. Add pseudo points

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_{ff} & K_{fu} \\ K_{uf} & K_{uu} \end{bmatrix} \right)$$



## 2. drop direct f-f dependencies

## 3. calibrate model

(e.g. using KL divergence, many choices)

$$\arg \min_{q(\mathbf{u}), \{q(\mathbf{f}_t | \mathbf{u})\}_{t=1}^T} \text{KL}(p(\mathbf{f}, \mathbf{u}) || q(\mathbf{u}) \prod_{t=1}^T q(\mathbf{f}_t | \mathbf{u})) \implies \begin{aligned} q(\mathbf{u}) &= p(\mathbf{u}) \\ q(\mathbf{f}_t | \mathbf{u}) &= p(\mathbf{f}_t | \mathbf{u}) \end{aligned}$$

equal to exact conditionals

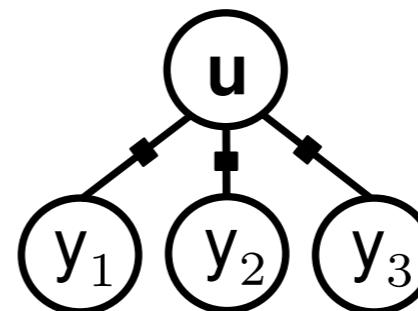
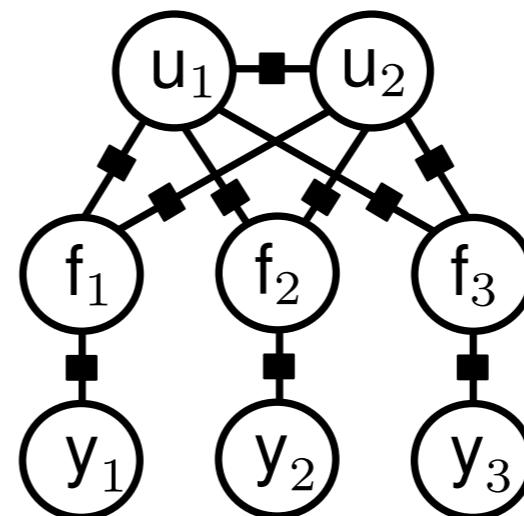
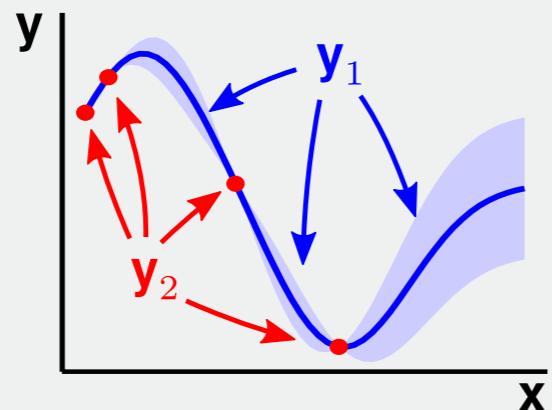
## 4. Replace p with simpler model

$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathbf{K}_{\mathbf{u}\mathbf{u}})$$

$$q(f_t | \mathbf{u}) = p(f_t | \mathbf{u})$$

How do we make predictions?

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \Sigma_{12} \Sigma_{22}^{-1} \mathbf{y}_2, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^\top)$$

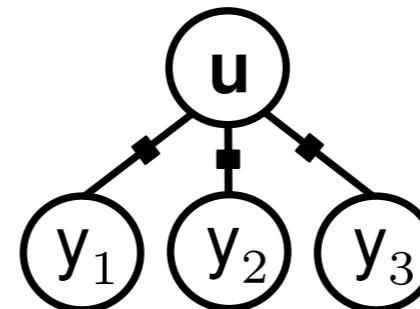
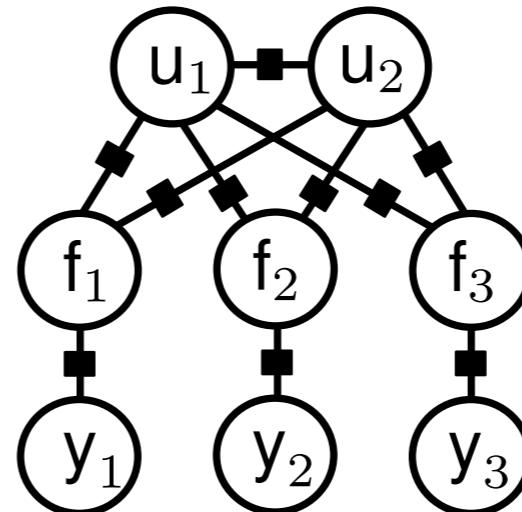
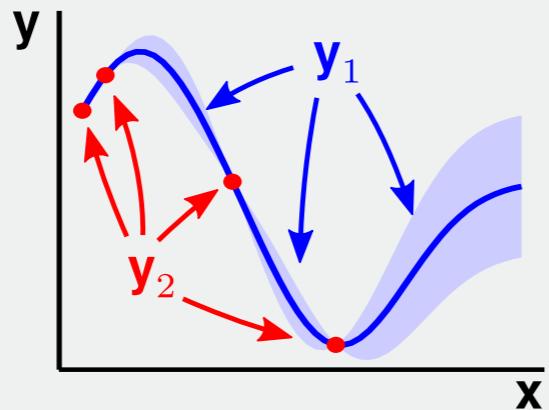


$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathbf{K}_{\mathbf{uu}})$$

$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

How do we make predictions?

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \Sigma_{12}\Sigma_{22}^{-1}\mathbf{y}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^\top)$$



$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

$$= \mathcal{N}(\mathbf{f}_t; \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{u}, \underbrace{\mathbf{K}_{\mathbf{f}_t \mathbf{f}_t} - \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{f}_t}}_{\mathbf{D}_{tt}})$$

$$q(\mathbf{y}_t | \mathbf{f}_t) = p(\mathbf{y}_t | \mathbf{f}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{f}_t, \sigma_y^2)$$

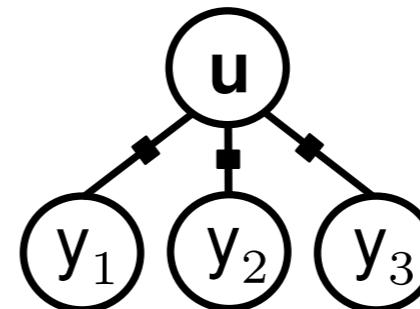
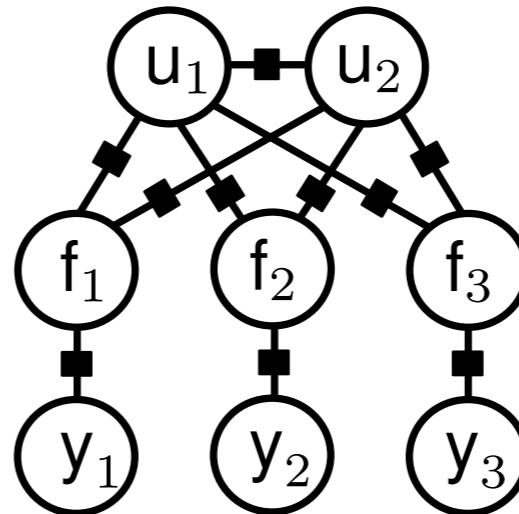
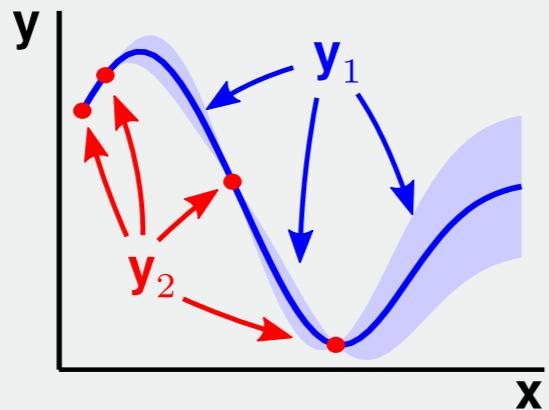
cost of computing likelihood is  $\mathcal{O}(TM^2)$

$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathbf{K}_{\mathbf{uu}})$$

$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

How do we make predictions?

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \Sigma_{12}\Sigma_{22}^{-1}\mathbf{y}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^\top)$$



$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

$$= \mathcal{N}(\mathbf{f}_t; \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{u}, \underbrace{\mathbf{K}_{\mathbf{f}_t \mathbf{f}_t} - \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{f}_t}}_{\mathbf{D}_{tt}})$$

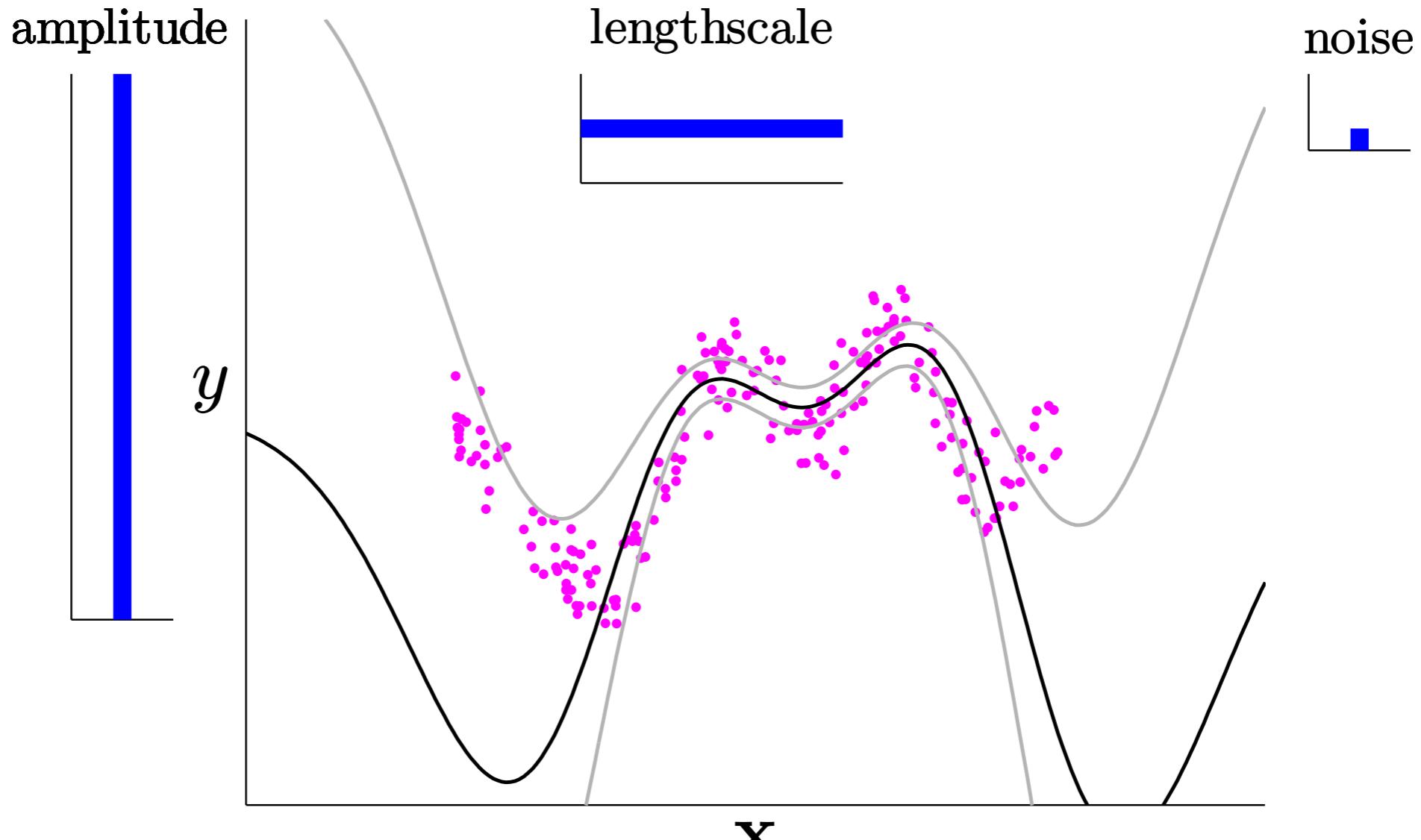
$$q(\mathbf{y}_t | \mathbf{f}_t) = p(\mathbf{y}_t | \mathbf{f}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{f}_t, \sigma_y^2)$$

cost of computing likelihood is  $\mathcal{O}(TM^2)$

$$p(\mathbf{y}_t | \theta) = \mathcal{N}(\mathbf{y}_t; \mathbf{0}, \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{f}_t}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{f}_t} + \mathbf{D} + \sigma_y^2 \mathbf{I})$$

cost of computing likelihood is  $\mathcal{O}(TM^2)$

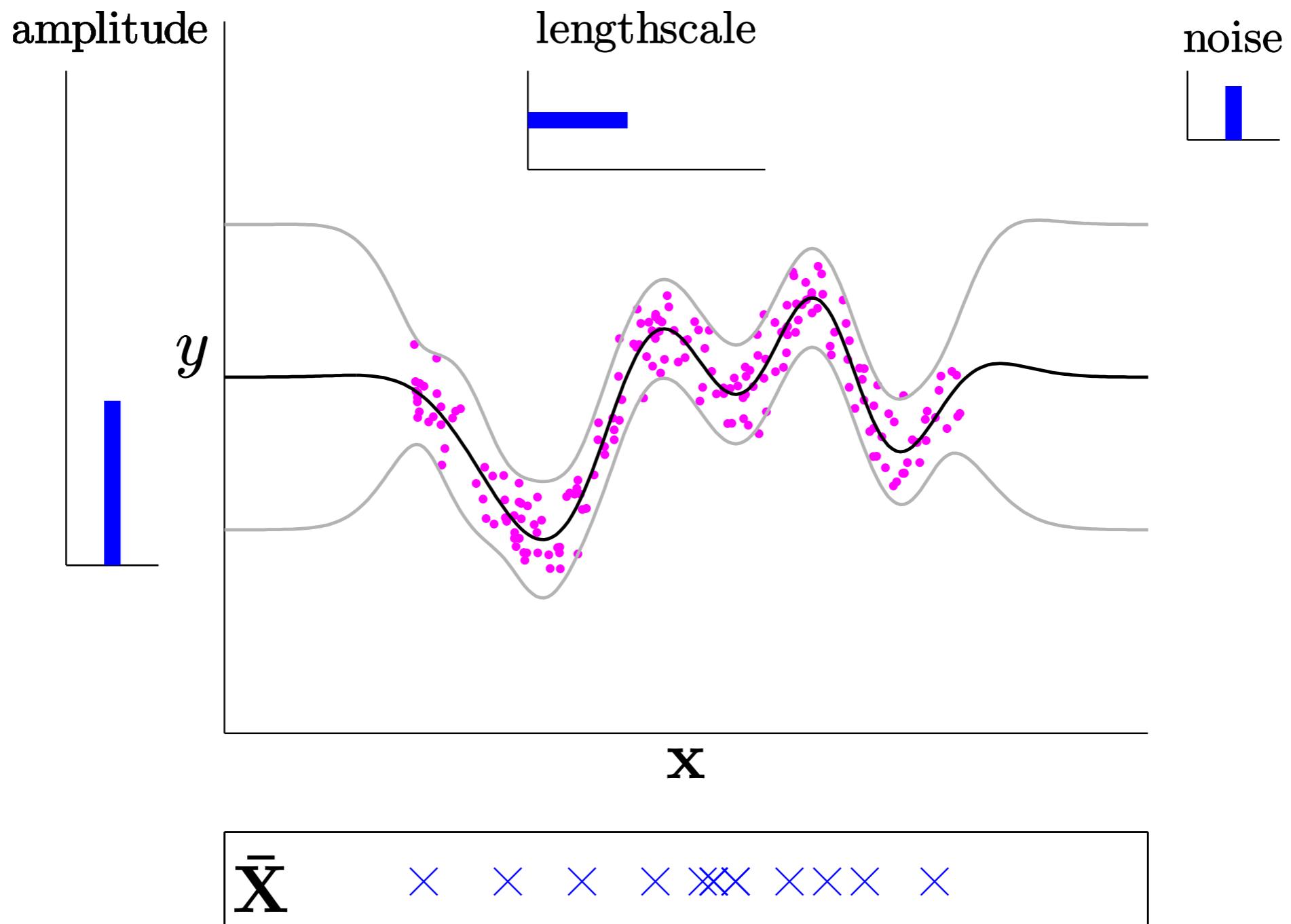
# FITC: Demo (Snelson)



Initialize adversarially:

amplitude and lengthscale too big  
noise too small  
pseudo-inputs bunched up

# FITC: Demo (Snelson)

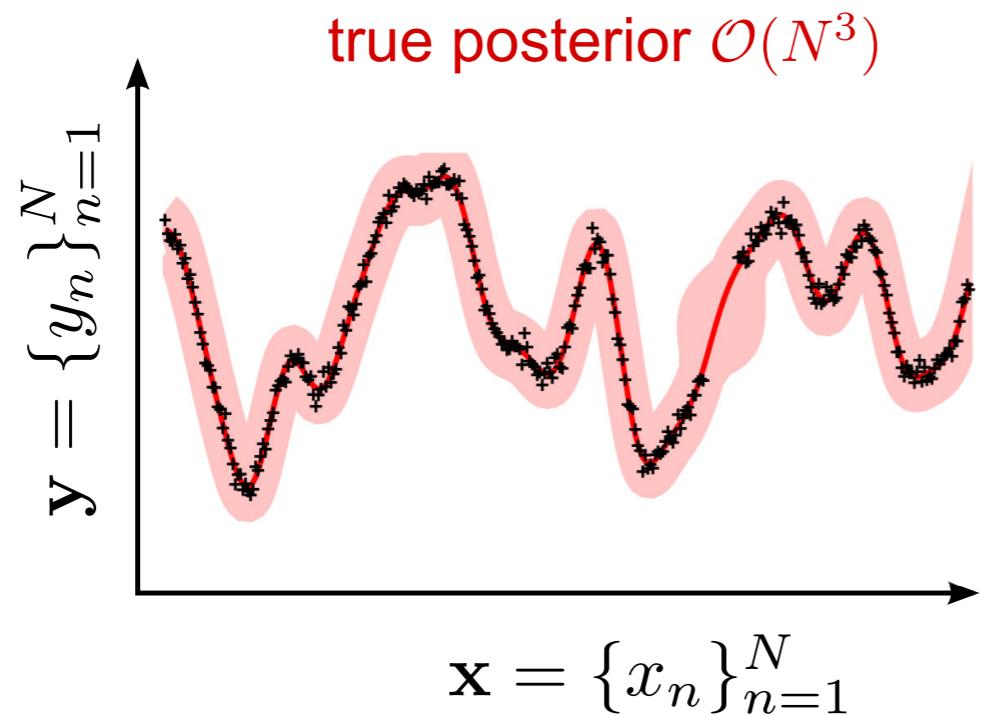


Pseudo-inputs and hyperparameters optimized

# EP pseudo-point approximation

---

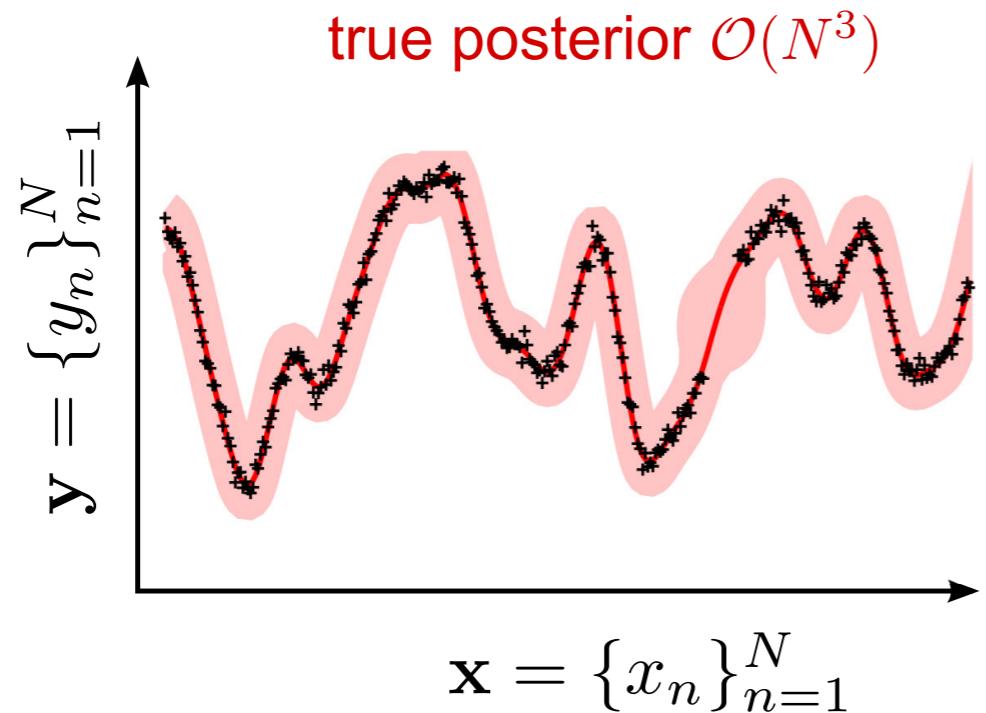
$$p^*(f) = p(f, \mathbf{y} | \mathbf{x}, \theta)$$



## EP pseudo-point approximation

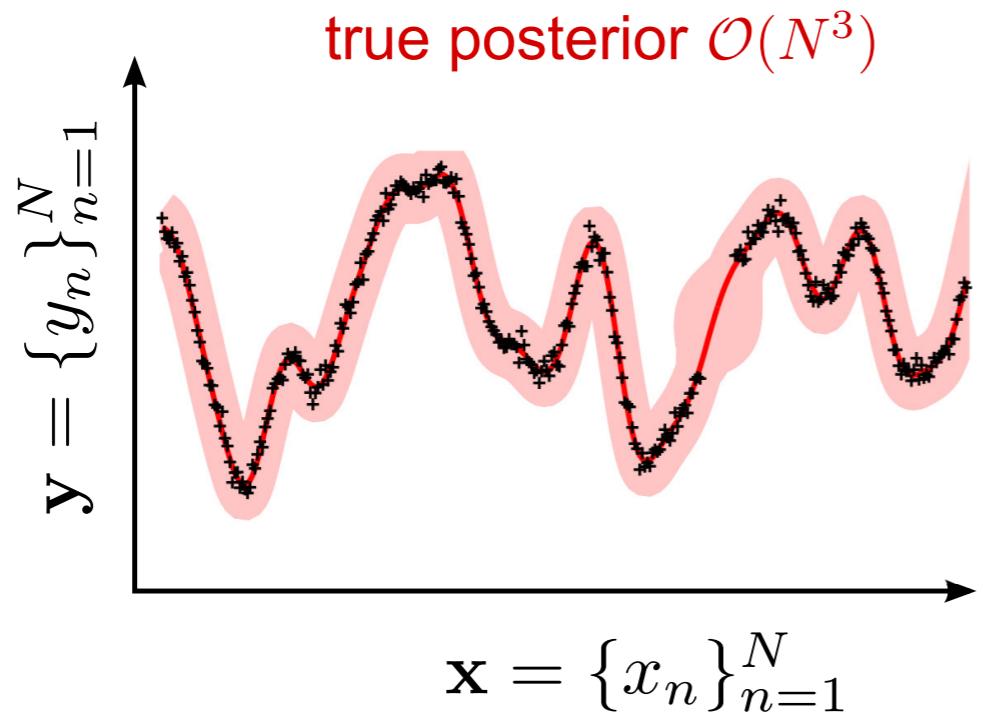
---

$$\begin{aligned} p^*(f) &= p(f, \mathbf{y} | \mathbf{x}, \theta) \\ &= p(f | \theta) \prod_{n=1}^N \underline{p(y_n | f, x_n, \theta)} \end{aligned}$$



## EP pseudo-point approximation

$$\begin{aligned} p^*(f) &= p(f, \mathbf{y} | \mathbf{x}, \theta) \\ &= p(f | \theta) \prod_{n=1}^N \underline{p(y_n | f, x_n, \theta)} \\ &= \underbrace{p(\mathbf{y} | \mathbf{x}, \theta)}_{\text{marginal likelihood}} \underbrace{p(f | \mathbf{y}, \mathbf{x}, \theta)}_{\text{posterior}} \end{aligned}$$



# EP pseudo-point approximation

$$p^*(f) = p(f, \mathbf{y} | \mathbf{x}, \theta)$$

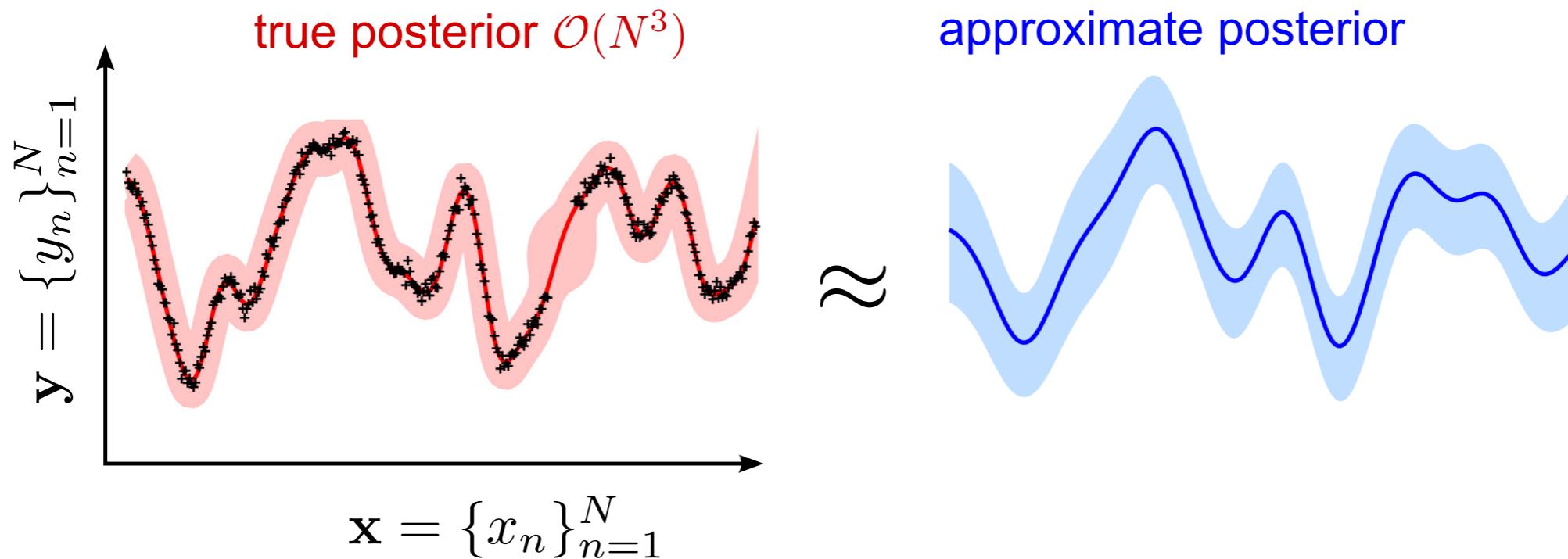
$$= p(f|\theta) \prod_{n=1}^N \underline{p(y_n|f, x_n, \theta)}$$

$$= \underline{p(\mathbf{y}|\mathbf{x}, \theta)} \underline{p(f|\mathbf{y}, \mathbf{x}, \theta)}$$

marginal  
likelihood

posterior

$$q^*(f) = p(f|\theta) \prod_{n=1}^N \underline{t_n(f)}$$

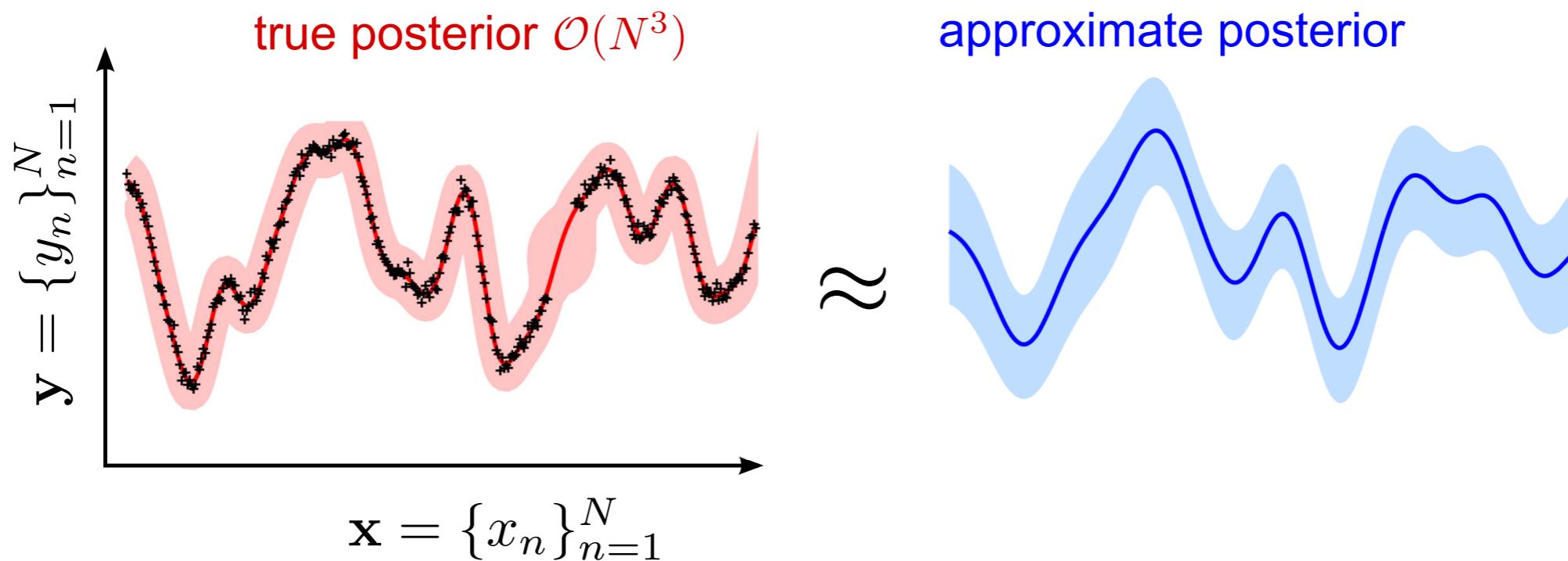


# EP pseudo-point approximation

$$p^*(f) = p(f, \mathbf{y} | \mathbf{x}, \theta)$$

$$\begin{aligned} &= p(f|\theta) \prod_{n=1}^N \underline{p(y_n|f, x_n, \theta)} \\ &= \underline{p(\mathbf{y}|\mathbf{x}, \theta)} \underline{p(f|\mathbf{y}, \mathbf{x}, \theta)} \\ &\quad \text{marginal likelihood} \qquad \text{posterior} \end{aligned}$$

$$\begin{aligned} q^*(f) &= p(f|\theta) \prod_{n=1}^N \underline{t_n(f)} \\ &= \underline{Z_{\text{EP}}} \underline{q(f)} \end{aligned}$$



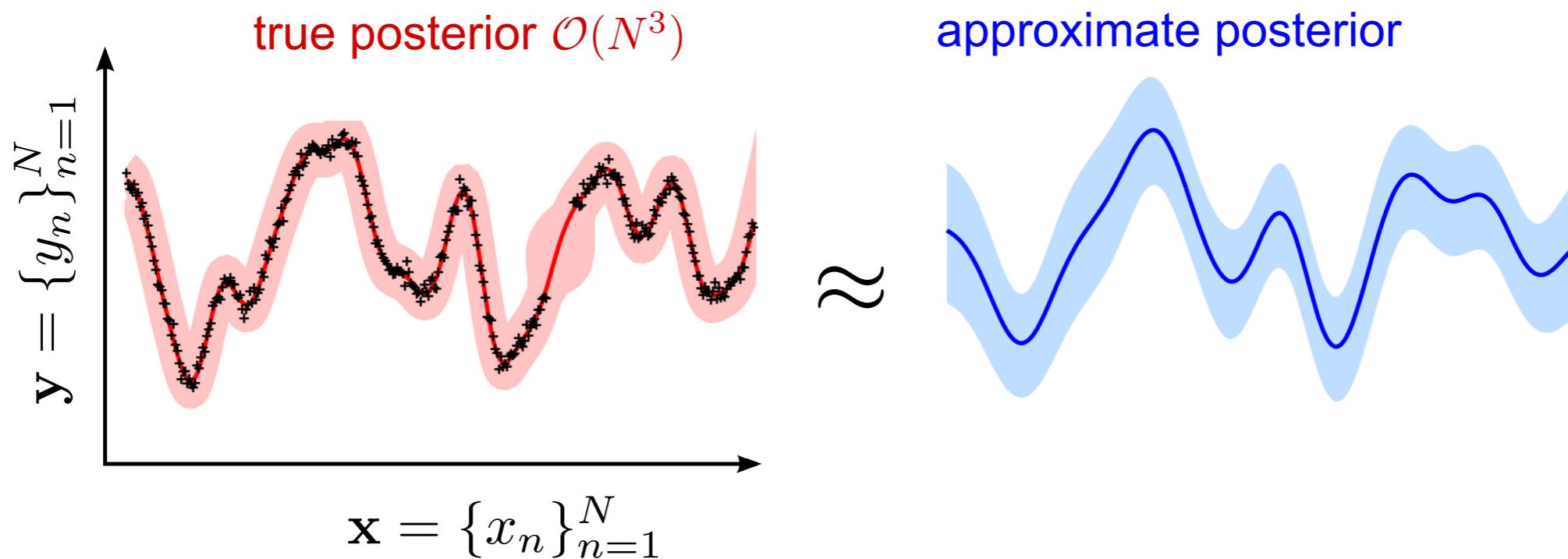
# EP pseudo-point approximation

$$\begin{aligned}
 p^*(f) &= p(f, \mathbf{y} | \mathbf{x}, \theta) \\
 &= p(f | \theta) \prod_{n=1}^N \underline{p(y_n | f, x_n, \theta)} \\
 &= \underline{p(\mathbf{y} | \mathbf{x}, \theta)} \underline{p(f | \mathbf{y}, \mathbf{x}, \theta)} \\
 &\quad \text{marginal likelihood} \qquad \text{posterior}
 \end{aligned}$$

$$\begin{aligned}
 q^*(f) &= p(f | \theta) \prod_{n=1}^N \underline{t_n(f)} \\
 &= \underline{Z_{\text{EP}}} \underline{q(f)}
 \end{aligned}$$

$$t_n(f) = \mathcal{N}(\mathbf{u}; \mu_n, \Sigma_n)$$

$$\dim(\mathbf{u}) = M \quad f = \{\mathbf{u}, f_{\neq \mathbf{u}}\}$$



# EP pseudo-point approximation

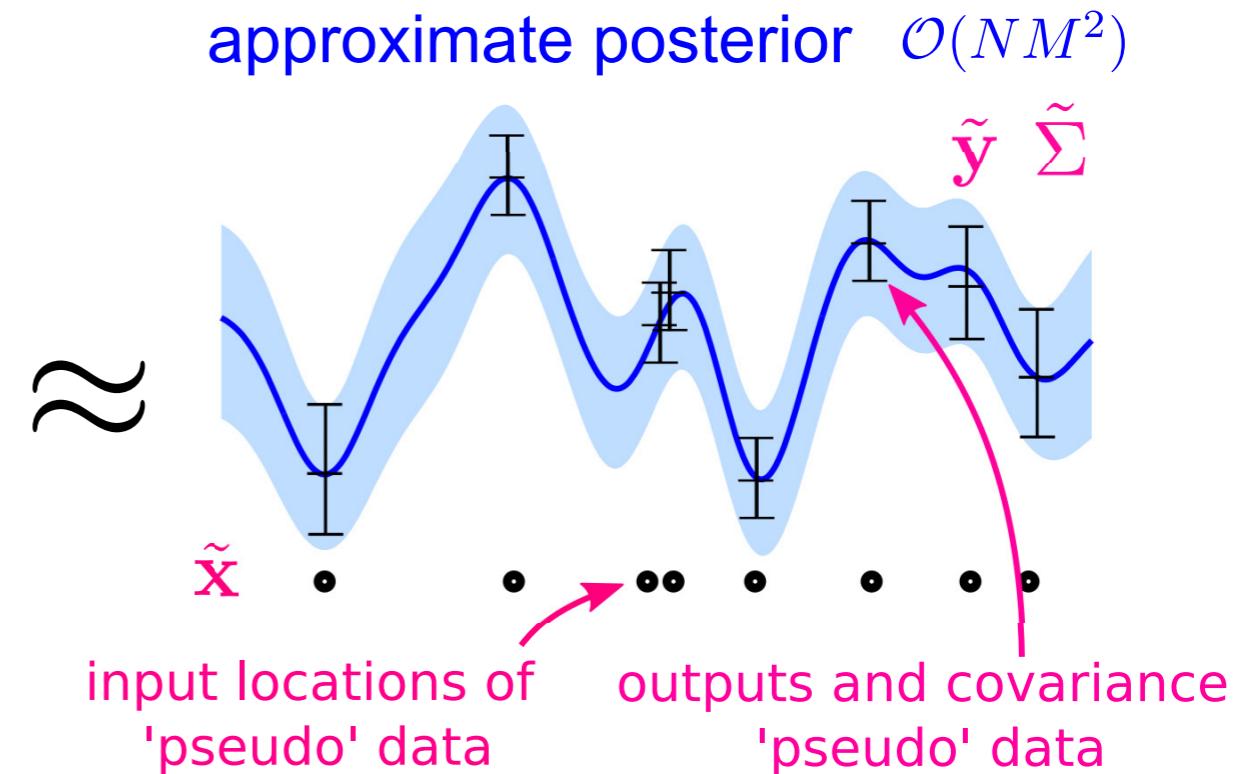
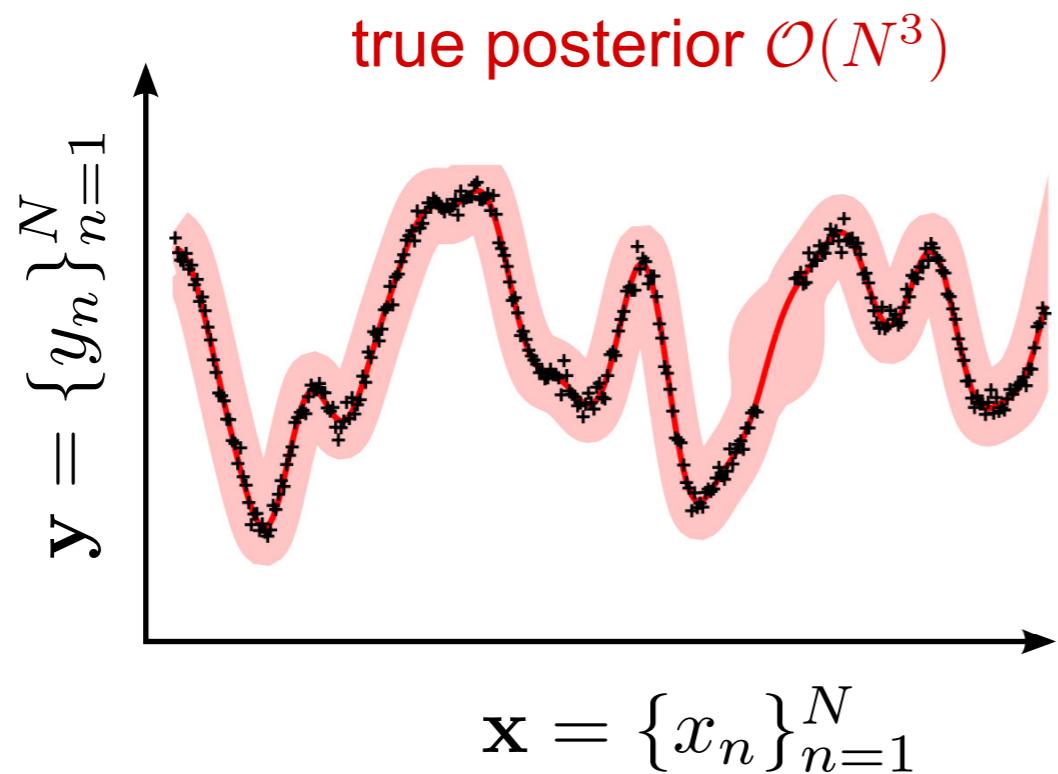
$$\begin{aligned}
 p^*(f) &= p(f, \mathbf{y} | \mathbf{x}, \theta) \\
 &= p(f | \theta) \prod_{n=1}^N \underline{p(y_n | f, x_n, \theta)} \\
 &= \underline{p(\mathbf{y} | \mathbf{x}, \theta)} \underline{p(f | \mathbf{y}, \mathbf{x}, \theta)} \\
 &\quad \text{marginal likelihood} \qquad \text{posterior}
 \end{aligned}$$

$$\begin{aligned}
 q^*(f) &= p(f | \theta) p(\tilde{\mathbf{y}} | \mathbf{u}, \tilde{\Sigma}) \\
 &= p(f | \theta) \prod_{n=1}^N \underline{t_n(f)} \\
 &= \underline{Z_{\text{EP}}} \underline{q(f)}
 \end{aligned}$$

exact joint  
of new GP  
regression  
model

$$t_n(f) = \mathcal{N}(\mathbf{u}; \mu_n, \Sigma_n)$$

$$\dim(\mathbf{u}) = M \quad f = \{\mathbf{u}, f_{\neq \mathbf{u}}\}$$



# EP algorithm

---

1. remove

$$q^{\setminus n}(f) = \frac{q^*(f)}{t_n(\mathbf{u})}$$

cavity



take out one  
pseudo-observation  
likelihood

# EP algorithm

---

1. remove

$$q^{\setminus n}(f) = \frac{q^*(f)}{t_n(\mathbf{u})}$$

cavity

take out one  
pseudo-observation  
likelihood

2. include

$$p_n^{\text{tilt}}(f) = q^{\setminus n}(f)p(y_n|f, x_n, \theta)$$

↑  
tilted

add in one  
true observation  
likelihood

# EP algorithm

---

1. remove

$$q^{\setminus n}(f) = \frac{q^*(f)}{t_n(\mathbf{u})}$$

cavity



take out one  
pseudo-observation  
likelihood

2. include

$$p_n^{\text{tilt}}(f) = q^{\setminus n}(f)p(y_n|f, x_n, \theta)$$

tilted



add in one  
true observation  
likelihood

3. project

$$q^*(f) = \operatorname{argmin}_{q^*(f)} \text{KL} [p_n^{\text{tilt}}(f) || q^*(f)]$$

project onto  
approximating  
family

KL between unnormalised  
stochastic processes



1. minimum: moments matched at pseudo-inputs  $\mathcal{O}(NM^2)$
2. Gaussian regression: matches moments everywhere

# EP algorithm

---

1. remove

$$q^{\setminus n}(f) = \frac{q^*(f)}{t_n(\mathbf{u})}$$

cavity



take out one  
pseudo-observation  
likelihood

2. include

$$p_n^{\text{tilt}}(f) = q^{\setminus n}(f)p(y_n|f, x_n, \theta)$$

tilted



add in one  
true observation  
likelihood

3. project

$$q^*(f) = \underset{q^*(f)}{\operatorname{argmin}} \text{KL} [p_n^{\text{tilt}}(f) || q^*(f)]$$

project onto  
approximating  
family

- 1. minimum: moments matched at pseudo-inputs  $\mathcal{O}(NM^2)$
- 2. Gaussian regression: matches moments everywhere

4. update

$$\begin{aligned} t_n(\mathbf{u}) &= \frac{q^*(f)}{q^{\setminus n}(f)} \\ &= z_n \mathcal{N}(\mathbf{K}_{f_n} \mathbf{u} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{u}; g_n, v_n) \end{aligned}$$

update  
pseudo-observation  
likelihood  
**rank 1**

KL between unnormalised  
stochastic processes



## **Further reading:**

### **Approximate inference in GPs:**

- Sparse Online Gaussian Processes, Csato and Opper, Neural Computation, 2002
- A Unifying View of Sparse Approximate Gaussian Process Regression, Quinonero-Candela and Rasmussen, JMLR, 2005
- Variational Learning of Inducing Variables in Sparse Gaussian Processes Titsias, AIStats, 2009
- On Sparse Variational Methods and the Kullback-Leibler Divergence between Stochastic Processes, Matthews et al., ICML 2016
- A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation, Bui et al., JMLR 2017
- Streaming Sparse Gaussian Process Approximations, Bui et al., NIPS 2017
- Efficient Deterministic Approximate Bayesian Inference for Gaussian Process Models , Bui, thesis, 2017

### **Deep Gaussian Processes:**

- Deep Gaussian Processes for Regression using Approximate Expectation Propagation, Bui et al., ICLR 2017
- Doubly Stochastic Variational Inference for Deep Gaussian Processes Salimbeni and Deisenroth, NIPS 2017

## **Gaussian Process Models**

- Gaussian Process Latent Variable Models for Visualisation of High Dimensional Data
- Lawrence, NIPS, 2005  
Local Distance Preservation in the GP-LVM through Back Constraints, Lawrence and
- Quinonero-Candela, ICML 2006

## **The Automatic Statistician**

- The Automatic Statistician, Ghahramani et al., ([website link](#))

---

## **GPs for Reinforcement Learning and Control**

- PILCO: A Model-Based and Data-Efficient Approach to Policy Search, Deisenroth and
- Rasmussen, ICML 2011

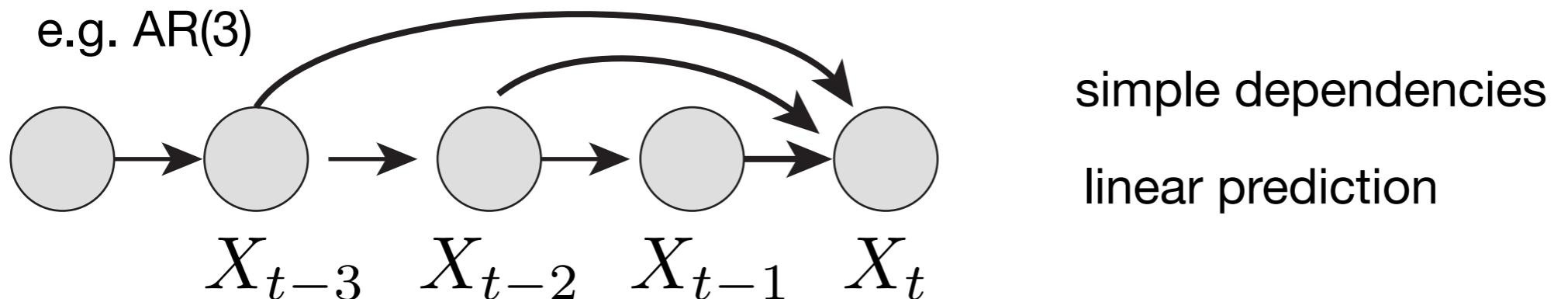
## **GPs and Neural Networks**

- Bayesian Learning for Neural Networks Neal, 1996
- Gaussian Process Behaviour in Wide Deep Neural Networks, Matthews et al., arXiv, 2018

**How do we use these to model time series?**

## How do we use these to model time series?

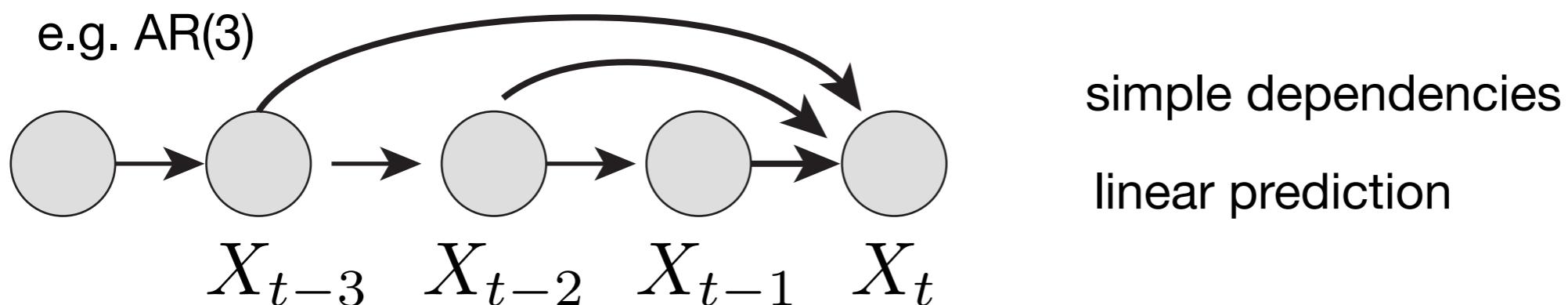
### ARIMA: directly model data statistics



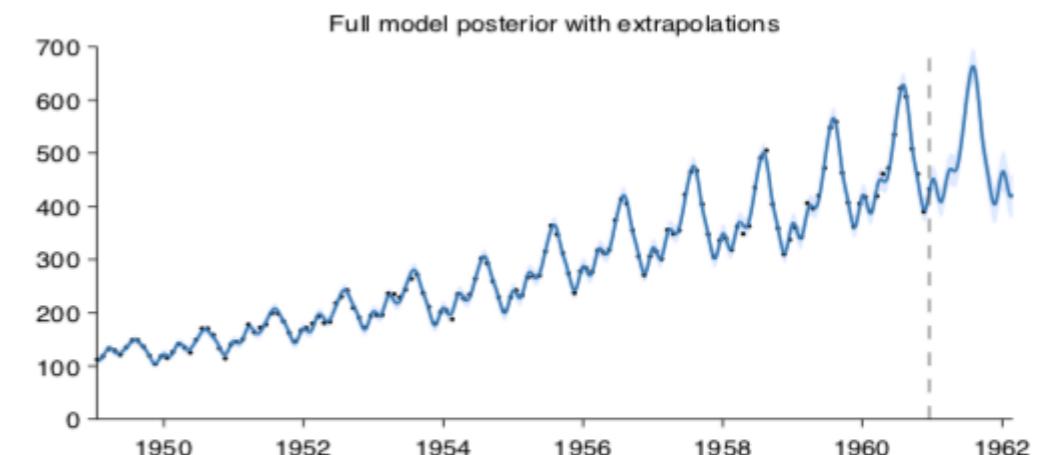
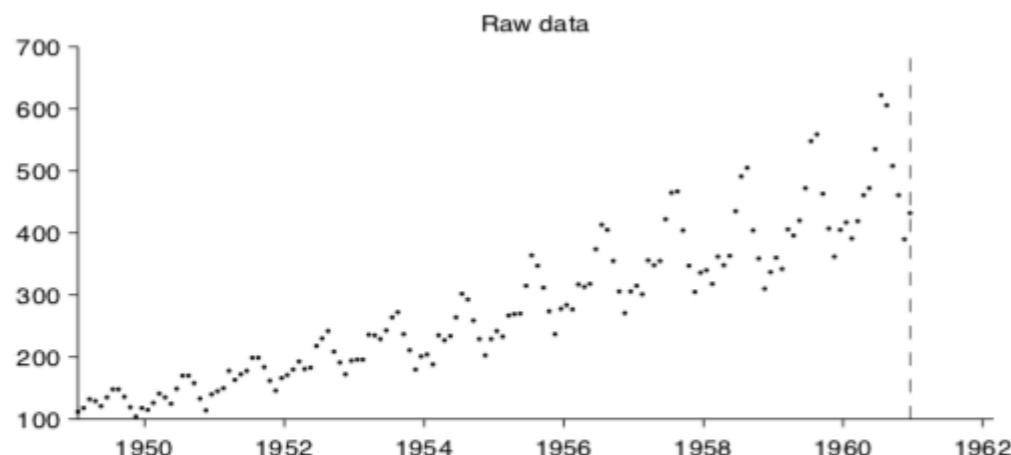
**GP: complex functional dependencies**

# How do we use these to model time series?

## ARIMA: directly model data statistics



## GP: complex functional dependencies



## Deep Kalman Filter (*Krishnan et al, 2016*)

Transition distribution

$$z_t | z_{t-1} \sim \mathcal{N}(F_\theta^\mu(z_{t-1}), F_\theta^\Sigma(z_{t-1}))$$

Emission distribution

$$x_t | z_t \sim \mathcal{N}(G_\theta^\mu(z_t), G_\theta^\Sigma(z_t))$$

Approximate posterior distribution

$$z_t | x_{1:T} \sim \mathcal{N}(R_\theta^\mu(x_{1:T}), R_\theta^\Sigma(x_{1:T}))$$

## Deep Kalman Filter (*Krishnan et al, 2016*)

Transition distribution

$$z_t | z_{t-1} \sim \mathcal{N}(F_\theta^\mu(z_{t-1}), F_\theta^\Sigma(z_{t-1}))$$

Emission distribution

$$x_t | z_t \sim \mathcal{N}(G_\theta^\mu(z_t), G_\theta^\Sigma(z_t))$$

Feedforward neural networks

Approximate posterior distribution

$$z_t | x_{1:T} \sim \mathcal{N}(R_\theta^\mu(x_{1:T}), R_\theta^\Sigma(x_{1:T}))$$

## Deep Kalman Filter (*Krishnan et al, 2016*)

Transition distribution

$$z_t | z_{t-1} \sim \mathcal{N}(F_\theta^\mu(z_{t-1}), F_\theta^\Sigma(z_{t-1}))$$

Emission distribution

$$x_t | z_t \sim \mathcal{N}(G_\theta^\mu(z_t), G_\theta^\Sigma(z_t))$$

Approximate posterior distribution

$$z_t | x_{1:T} \sim \mathcal{N}(R_\theta^\mu(x_{1:T}), R_\theta^\Sigma(x_{1:T}))$$

Feedforward neural networks

Bidirectional recurrent networks

## Deep Kalman Filter (*Krishnan et al, 2016*)

Transition distribution

$$z_t | z_{t-1} \sim \mathcal{N}(F_\theta^\mu(z_{t-1}), F_\theta^\Sigma(z_{t-1}))$$

Emission distribution

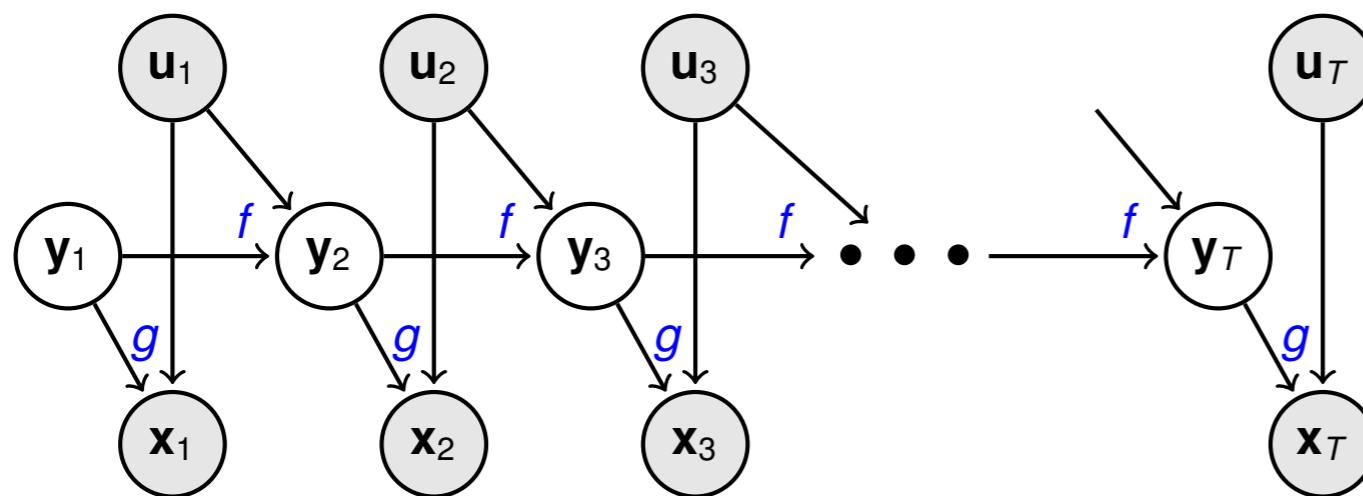
$$x_t | z_t \sim \mathcal{N}(G_\theta^\mu(z_t), G_\theta^\Sigma(z_t))$$

Approximate posterior distribution

$$z_t | x_{1:T} \sim \mathcal{N}(R_\theta^\mu(x_{1:T}), R_\theta^\Sigma(x_{1:T}))$$

Feedforward neural networks

Bidirectional recurrent networks



$$\begin{aligned}\mathbf{y}_{t+1} &= f(\mathbf{y}_t, \mathbf{u}_t) + \mathbf{w}_t \\ \mathbf{x}_t &= g(\mathbf{y}_t, \mathbf{u}_t) + \mathbf{v}_t\end{aligned}$$

$\mathbf{w}_t, \mathbf{v}_t$  usually assumed Gaussian.

# Deep Kalman Filter (*Krishnan et al, 2016*)

Transition distribution

$$z_t | z_{t-1} \sim \mathcal{N}(F_\theta^\mu(z_{t-1}), F_\theta^\Sigma(z_{t-1}))$$

Emission distribution

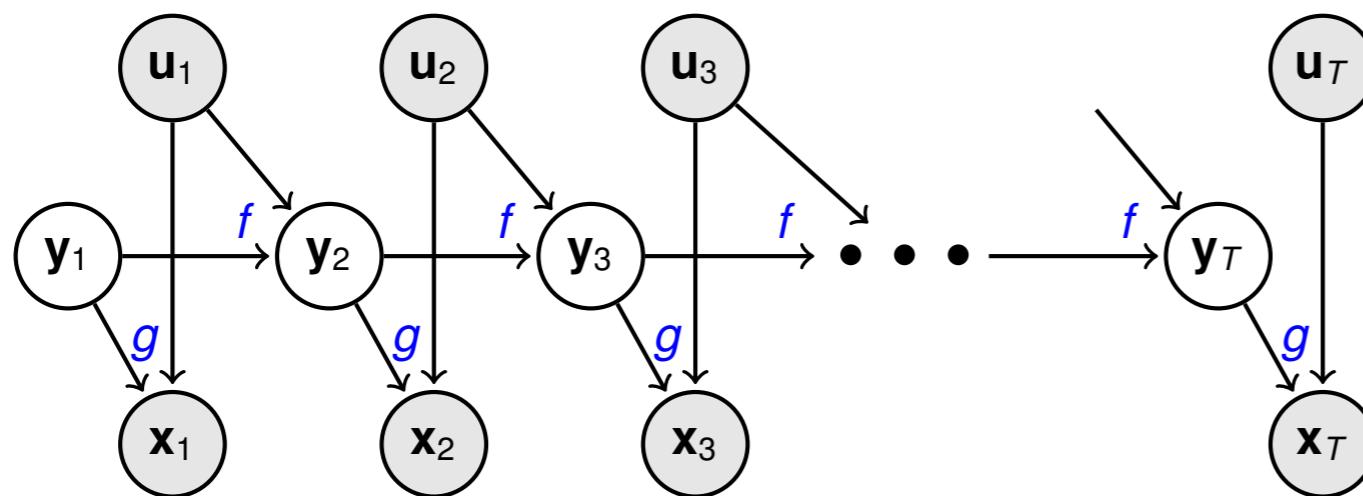
$$x_t | z_t \sim \mathcal{N}(G_\theta^\mu(z_t), G_\theta^\Sigma(z_t))$$

Approximate posterior distribution

$$z_t | x_{1:T} \sim \mathcal{N}(R_\theta^\mu(x_{1:T}), R_\theta^\Sigma(x_{1:T}))$$

Feedforward neural networks

Bidirectional recurrent networks

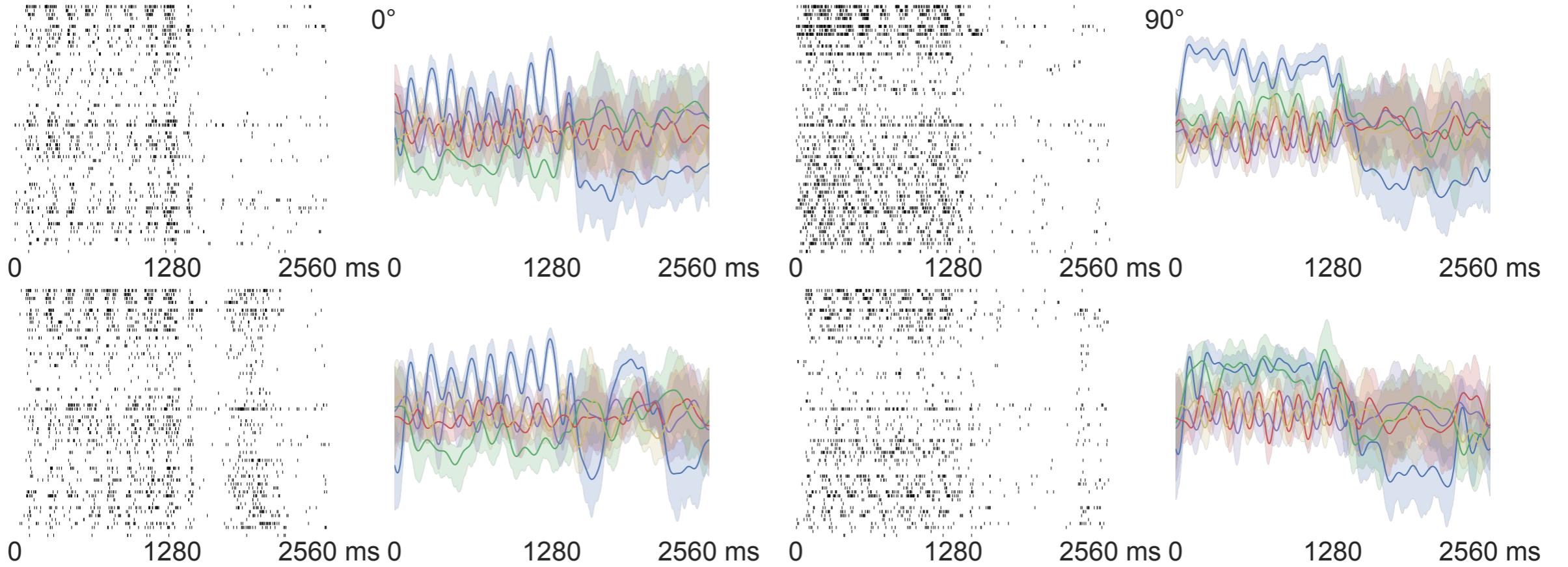


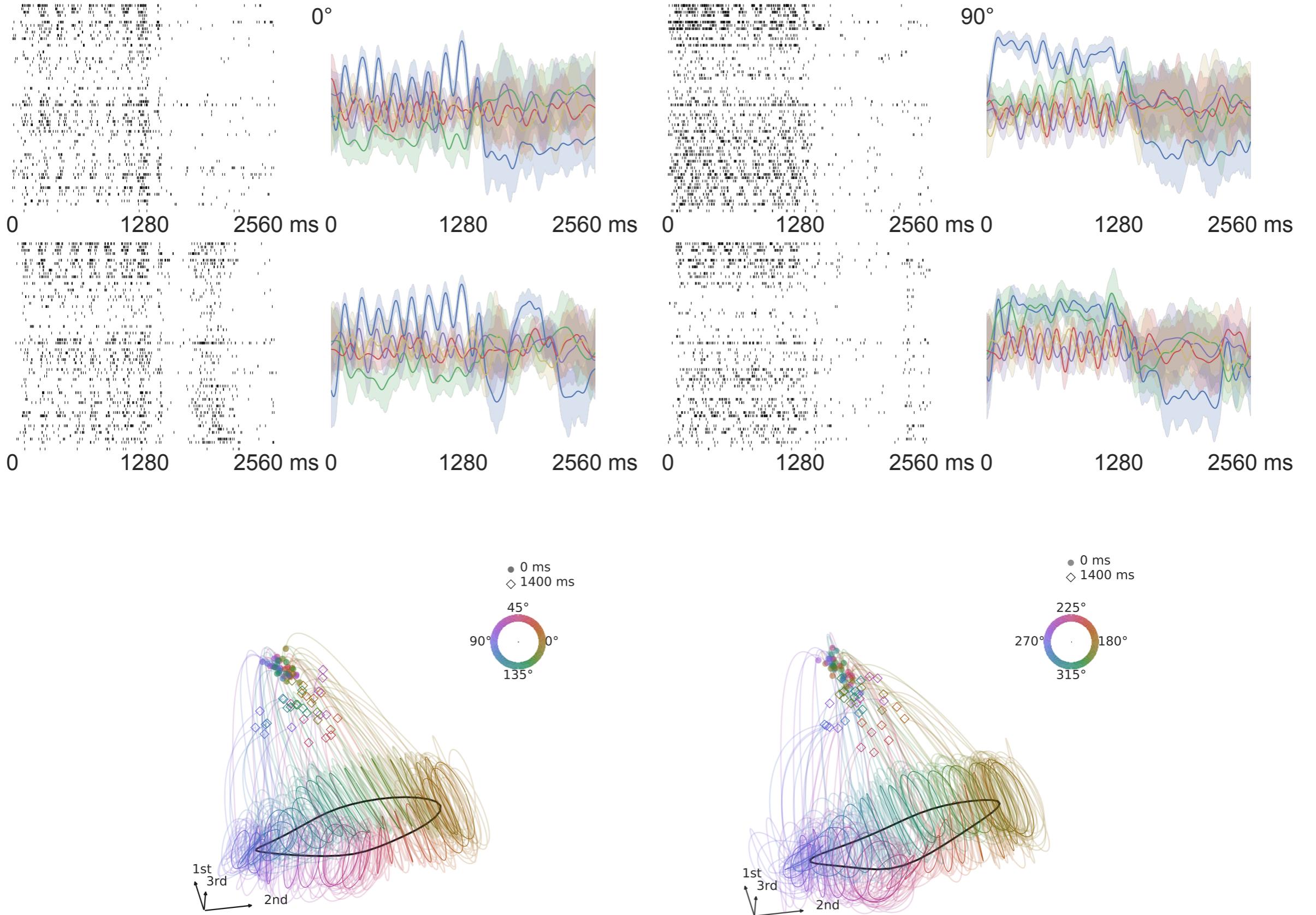
$$\mathbf{y}_{t+1} = f(\mathbf{y}_t, \mathbf{u}_t) + \mathbf{w}_t$$

$$\mathbf{x}_t = g(\mathbf{y}_t, \mathbf{u}_t) + \mathbf{v}_t$$

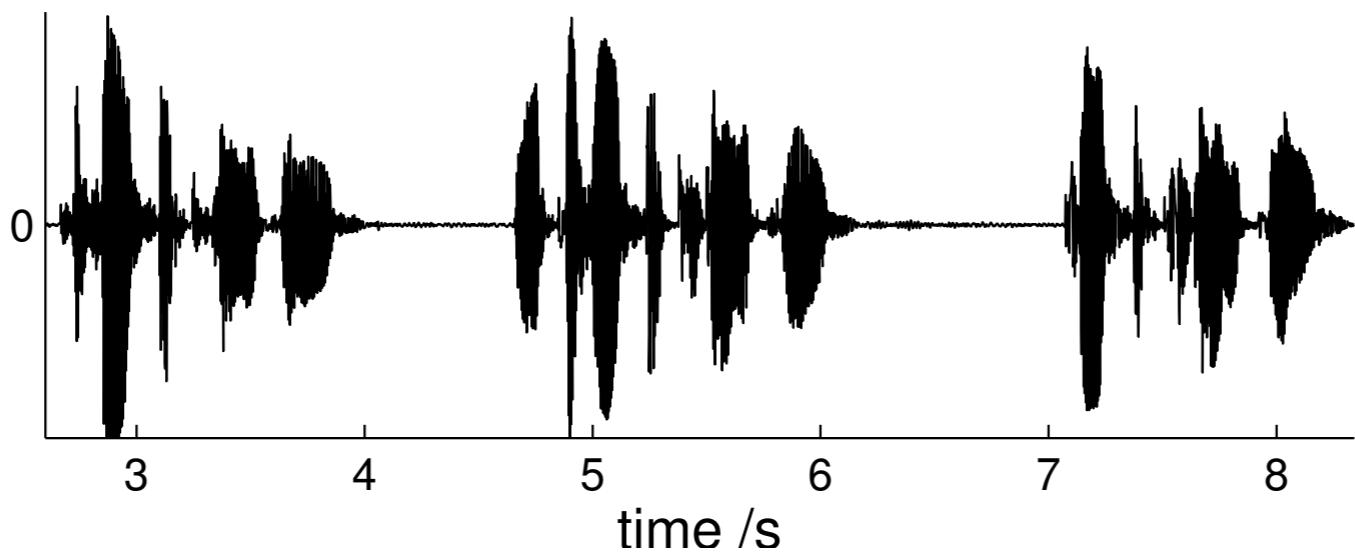
$\mathbf{w}_t, \mathbf{v}_t$  usually assumed Gaussian.

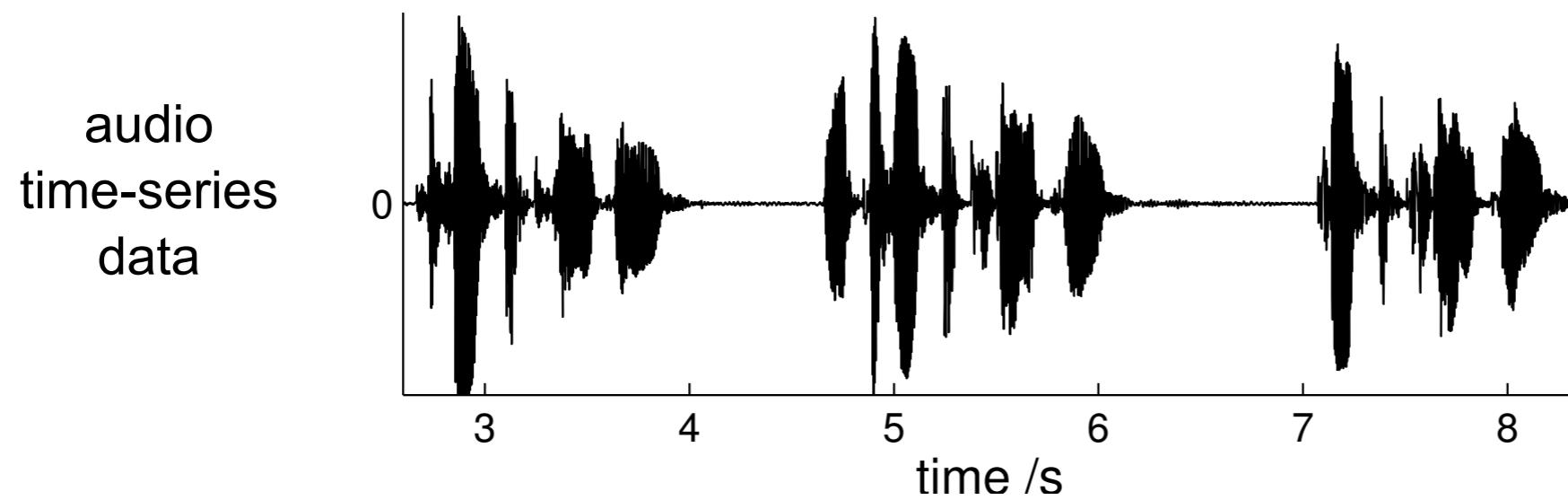
GP priors for f,g





audio  
time-series  
data





Generative model

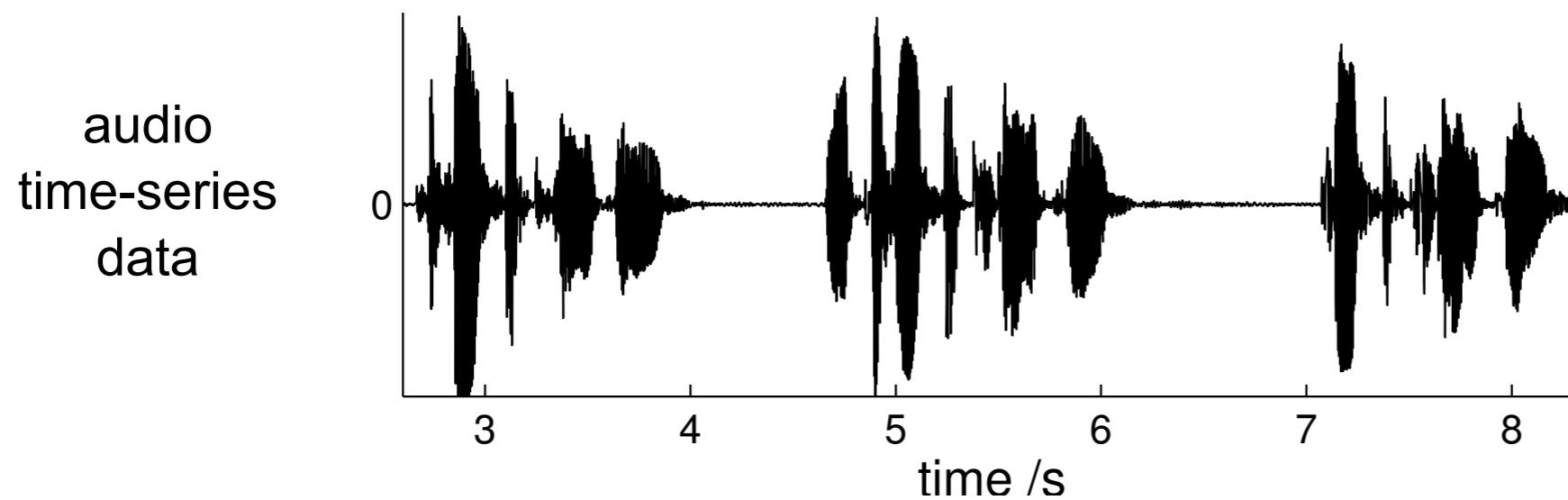
$$y(t) = f(t) + \epsilon \sigma_y \quad (\text{independent Gaussian noise})$$

$$p(\epsilon) = \mathcal{N}(0, 1)$$

place GP prior over the non-linear function

$$p(f(t)|\theta) = \mathcal{GP}(0, K(t, t'))$$

$$K(t, t') = \sigma^2 \cos(\omega(t - t')) \exp\left(-\frac{1}{2l^2}(t - t')^2\right)$$



Generative model

$$y(t) = f(t) + \epsilon\sigma_y \quad (\text{independent Gaussian noise})$$

$$p(\epsilon) = \mathcal{N}(0, 1)$$

place GP prior over the non-linear function

$$p(f(t)|\theta) = \mathcal{GP}(0, K(t, t'))$$

$$K(t, t') = \sigma^2 \cos(\omega(t - t')) \exp\left(-\frac{1}{2l^2}(t - t')^2\right)$$

**Also: model envelopes  
in spectral domain...**