

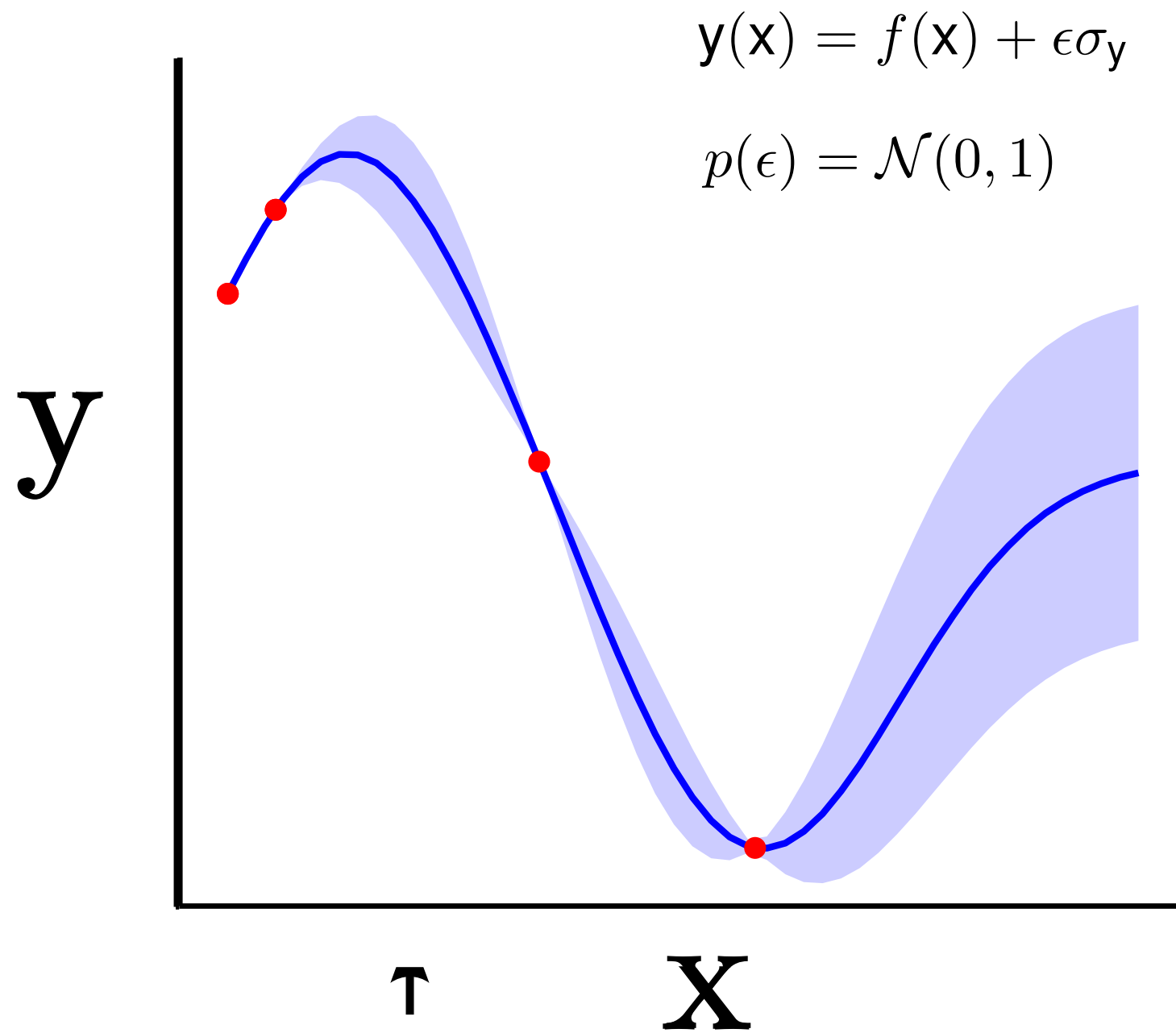
DS-GA 1810.001 Modeling time series data

L11. (From last time) Fast GP approximations

Instructor: Cristina Savin
NYU, CNS & CDS

Quick recap

1. Motivation: nonlinear regression

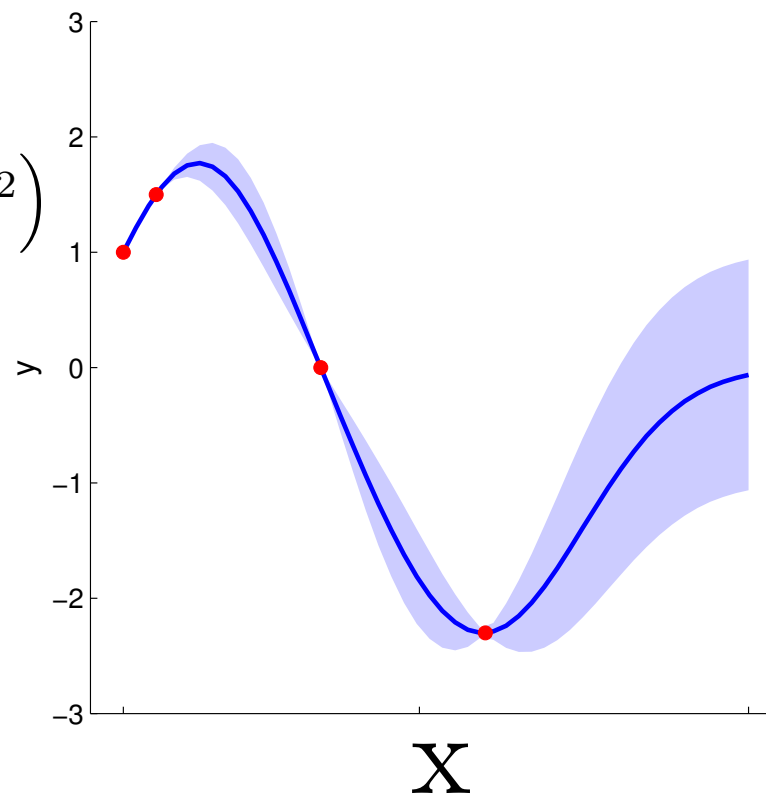
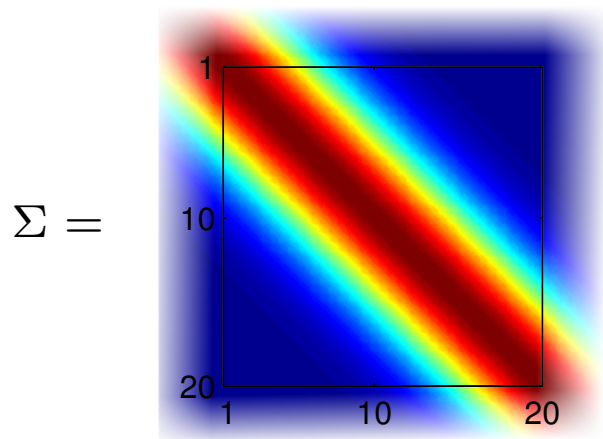


Quick recap

2.prior over functions: GP

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

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



GP: generalization of multivariate gaussian

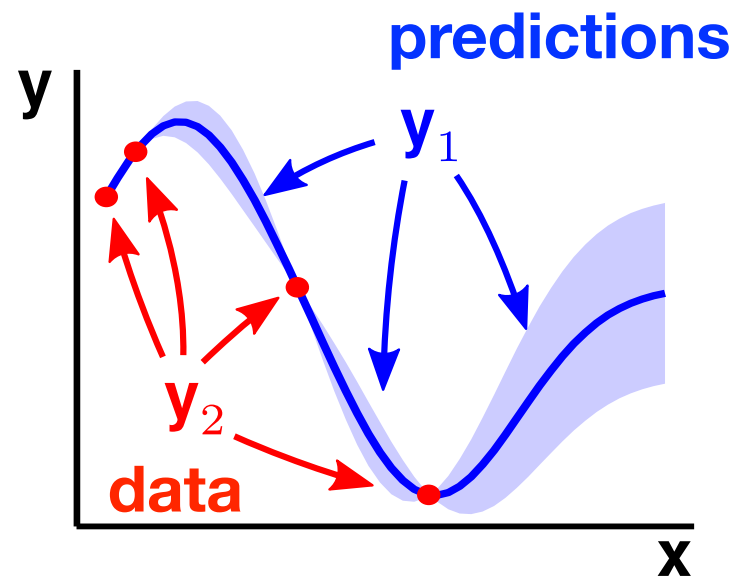
Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

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

Any function has nonzero probability, preference for certain structure

Quick recap

3. GP inference



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

Jointly gaussian:

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



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

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

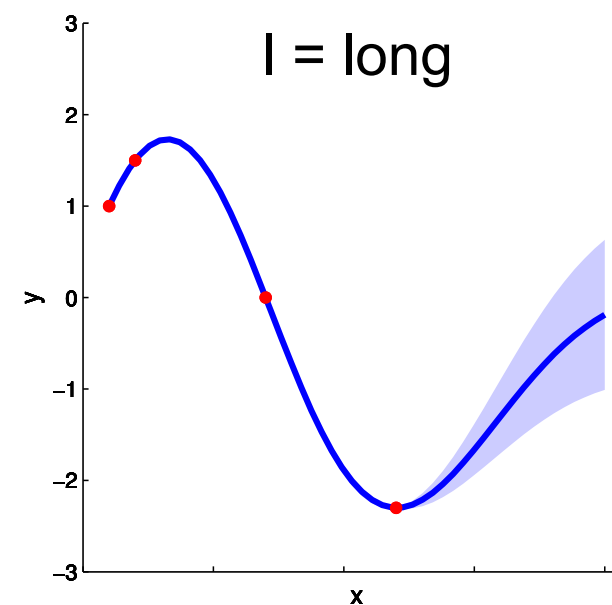
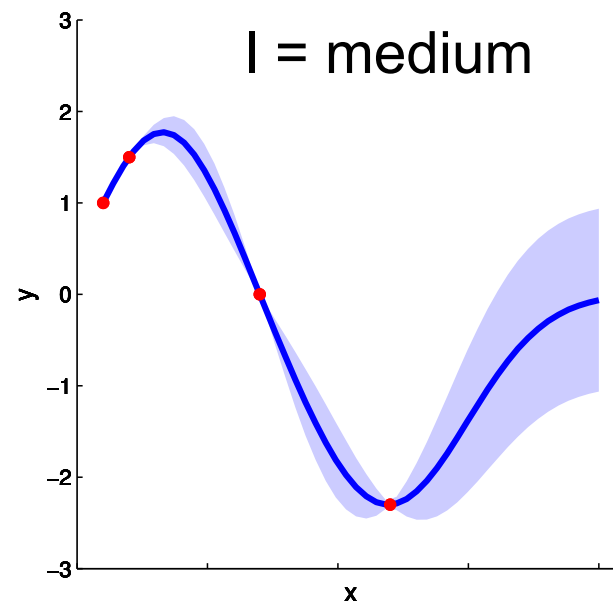
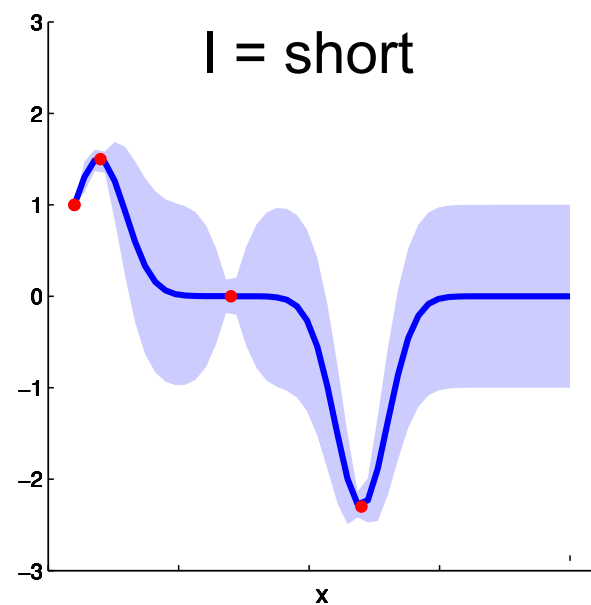
$$\Sigma_{\mathbf{y}_1 | \mathbf{y}_2} = \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^\top$$

**Computational
bottleneck!!!**

Quick recap

4. GP hyperparameters

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



Hyperparameters
significantly
influence outcome

Learn values from data!

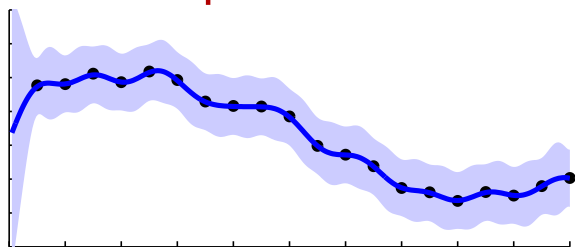
A. Maximum likelihood fit:

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

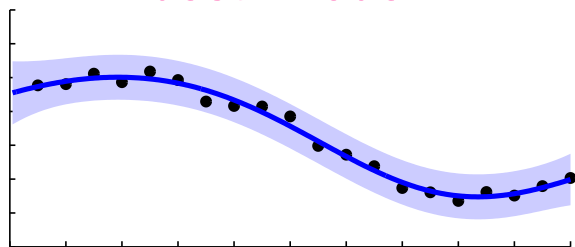
B. Bayesian:

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

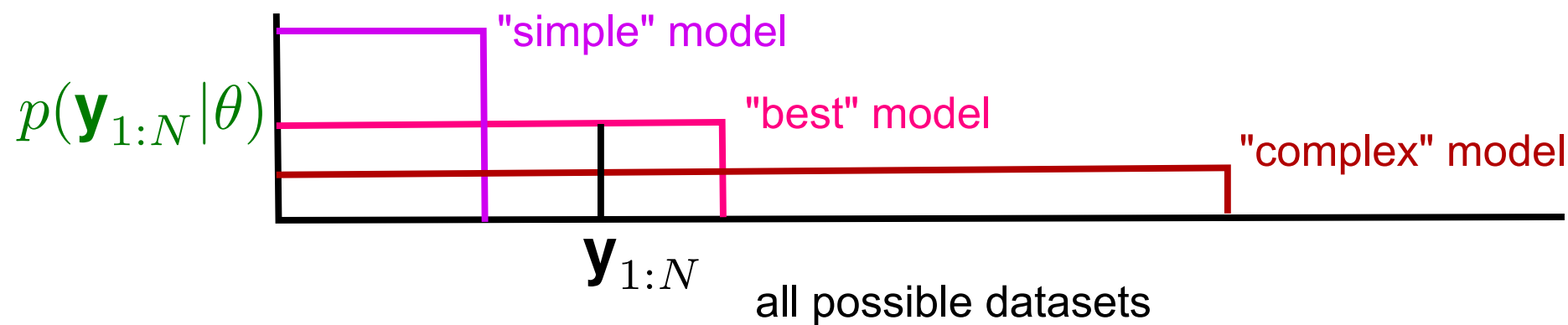
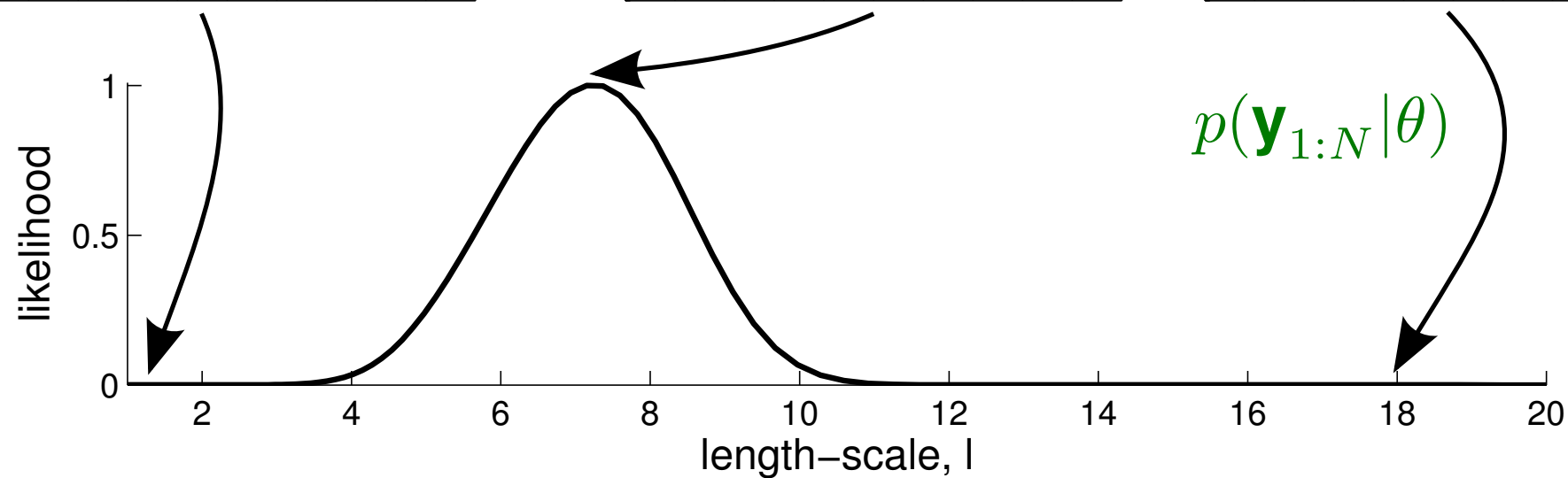
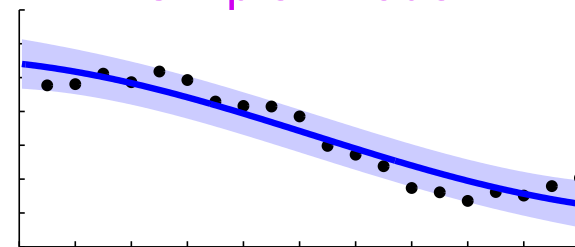
fits every training point
"complex" model



"best" model

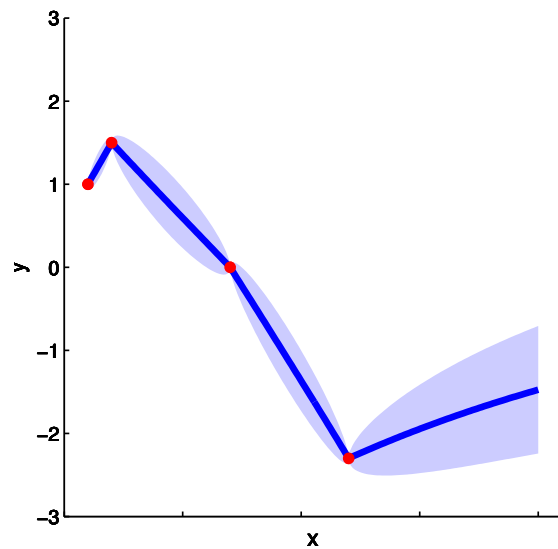


straight line-like
"simple" model

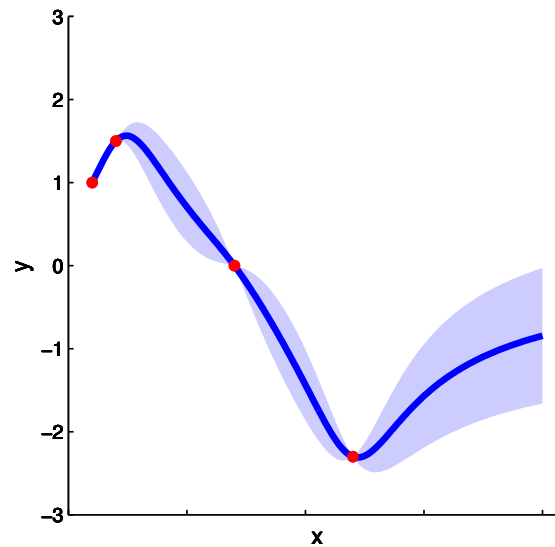


Quick recap 5. GP kernels

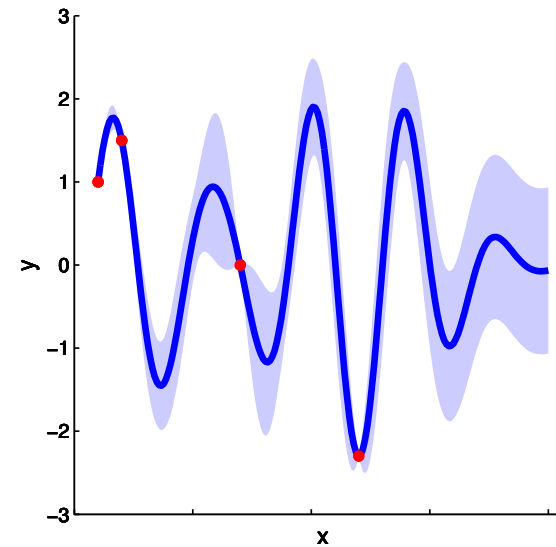
OU



RQ



periodic



Choice of
covariance function
influence outcome

$$p(M|\mathbf{y}_{1:N}) = \frac{\overset{\text{Marginal likelihood}}{p(\mathbf{y}_{1:N}|M)} \overset{\text{Prior over models}}{p(M)}}{\sum_{M'} p(\mathbf{y}_{1:N}|M') p(M')}$$

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

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)

A Brief History of Gaussian Process Approximations

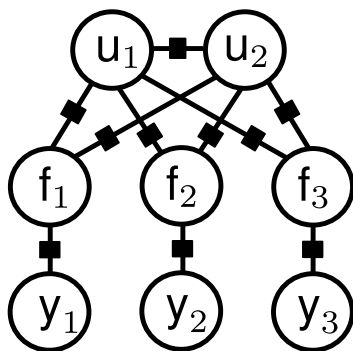
approximate generative model
exact inference

methods employing
pseudo-data

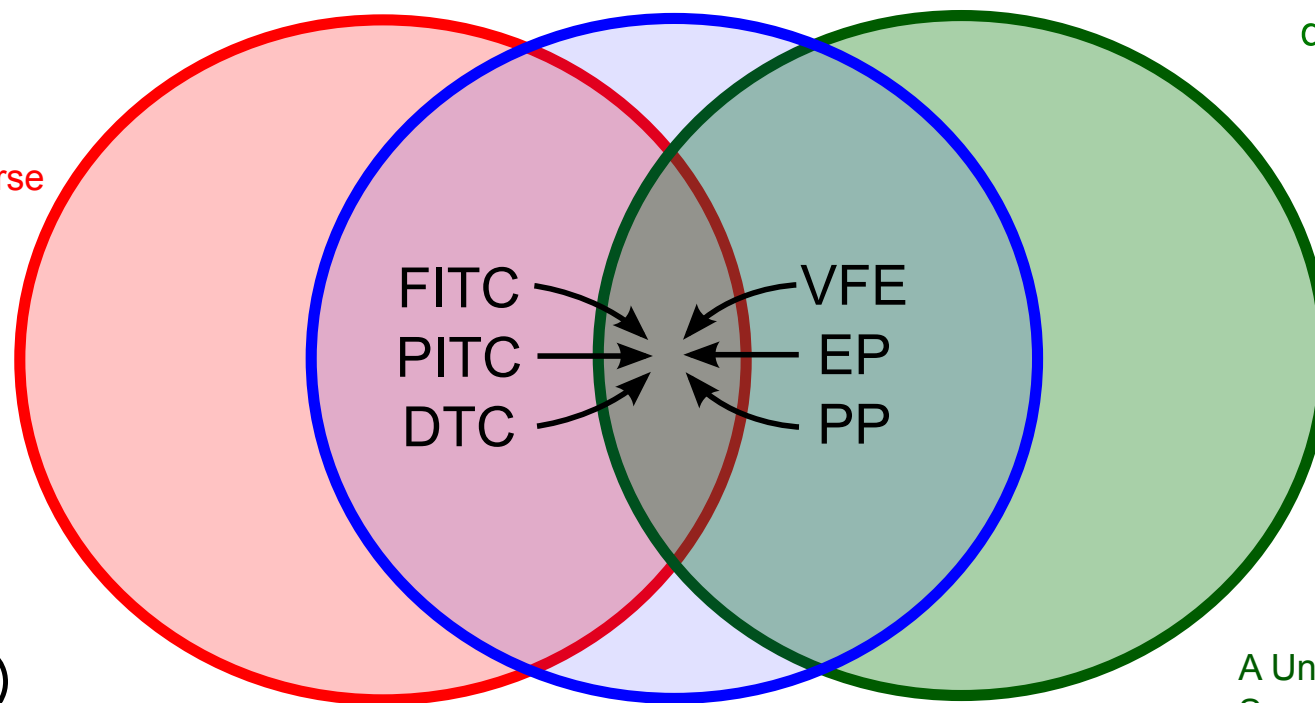
exact generative model
approximate inference

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

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



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



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"

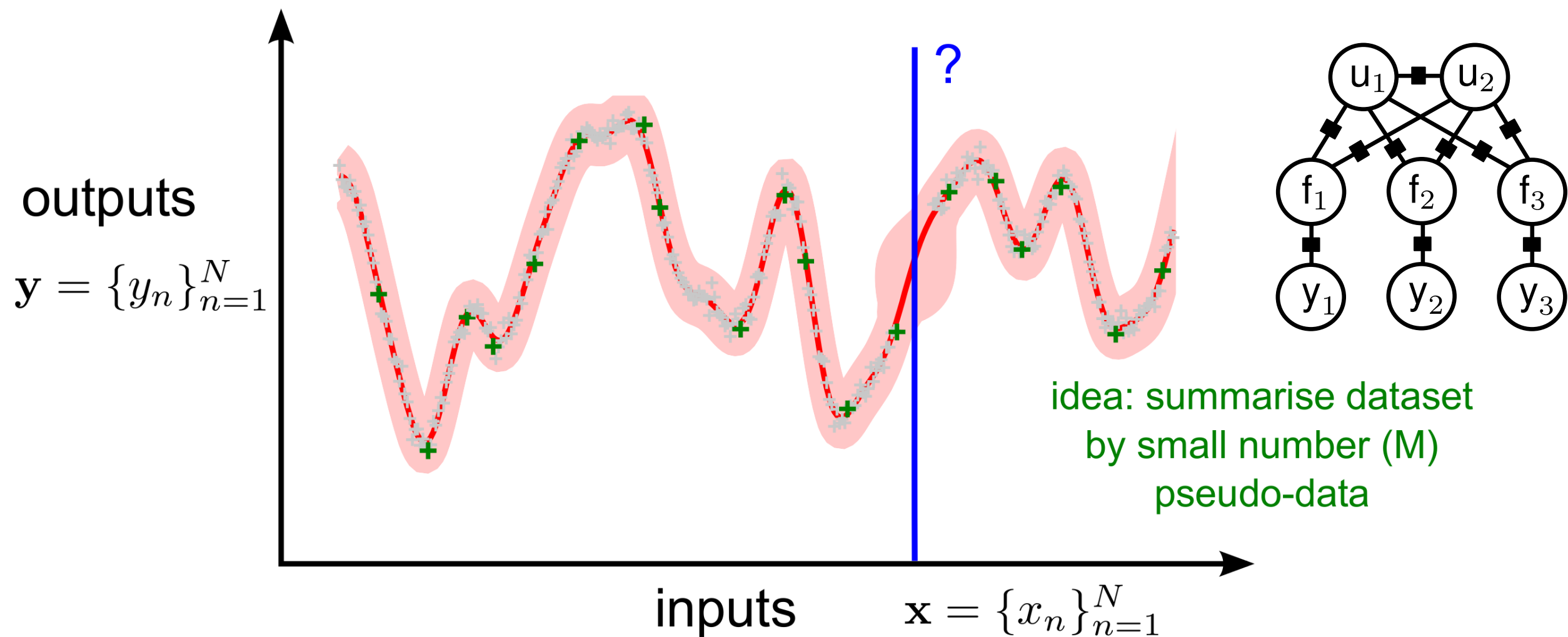
EP: Csato and Opper 2002 / Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

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

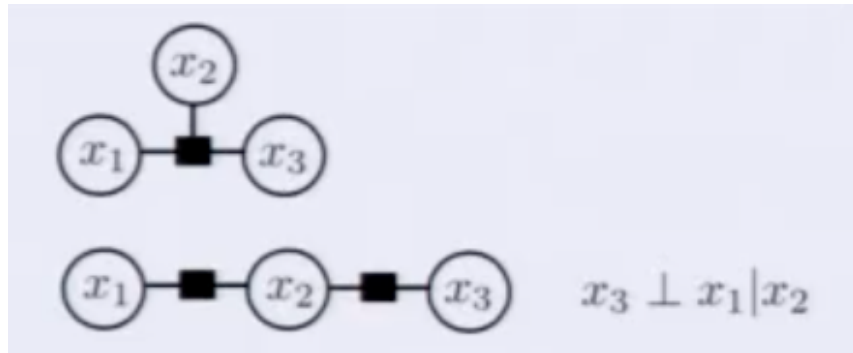
DTC / PP: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

$$p(f|\theta) = \mathcal{GP}(f; 0, K_\theta) \xrightarrow[\text{intractabilities, computational } \mathcal{O}(N^3), \text{ analytic}]{\text{inference \& learning}} p(f|\mathbf{y}, \mathbf{x}, \theta)$$

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



An interlude: dependencies in multivariate gaussians and factor graphs

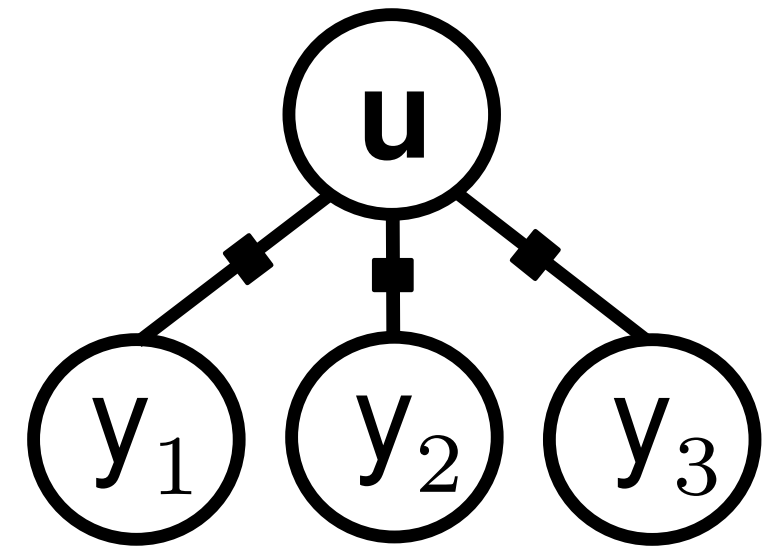


$$\Sigma^{-1} = \begin{bmatrix} 1.5 & -1/2 & -1/2 & 0 \\ -1/2 & 1 & 0 & 0 \\ -1/2 & 0 & 5/4 & -1/2 \\ 0 & 0 & -1/2 & 1 \end{bmatrix}$$

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

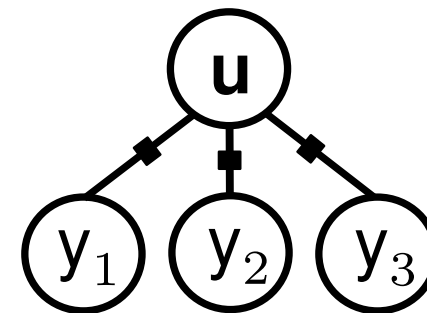
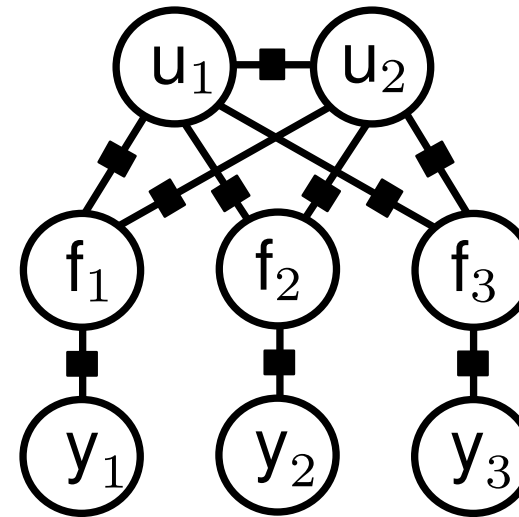
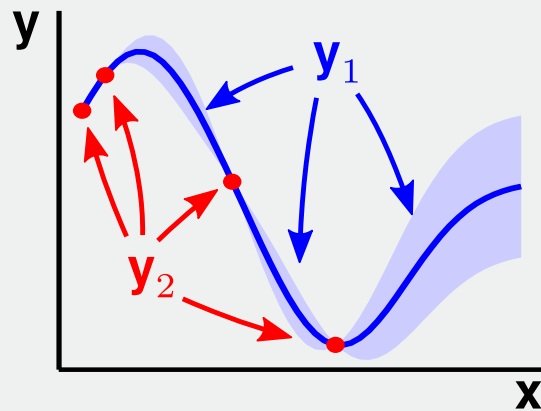
4. Replace p with simpler model

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

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

How do we make predictions?

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



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

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

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

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

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

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