

EE613 - Machine Learning for Engineers

<https://moodle.epfl.ch/course/view.php?id=16819>

NONLINEAR REGRESSION

Sylvain Calinon

Robot Learning and Interaction Group

Idiap Research Institute

Nov 9, 2023

EE613 schedule

Thu. 21.09.2023	(C) 1. ML introduction
Thu. 28.09.2023	(C) 2. Bayesian 1 (C) 3. Bayesian 2
Thu. 12.10.2023	(C) 4. Hidden Markov Models
Thu. 19.10.2023	(C) 5. Dimensionality reduction
Thu. 26.10.2023	(C) 6. Decision trees
Thu. 02.11.2023	(C) 7. Linear regression
Thu. 09.11.2023	(C) 8. Nonlinear regression
Thu. 16.11.2023	(C) 9. Kernel Methods - SVM
Thu. 23.11.2023	(C) 10. Tensor factorization
Thu. 30.11.2023	(C) 11. Deep learning 1
Thu. 07.12.2023	(C) 12. Deep learning 2
Thu. 14.12.2023	(C) 13. Deep learning 3
Thu. 21.12.2023	(C) 14. Deep learning 4

Properties of multivariate Gaussian distributions:

- Product of Gaussians
- Linear transformation and combination
- Conditional distribution

Three nonlinear regression models:

- Locally weighted regression (LWR)
- Gaussian mixture regression (GMR)
- Gaussian process regression (GPR)

**Modeling possible
co-variations**
(a.k.a. aleatoric uncertainty)

**Modeling uncertainty
of the estimate**
(a.k.a. epistemic uncertainty)

Multivariate Gaussian distribution

Univariate Gaussian distribution:

$$\mathcal{N}(\mu, \sigma^2) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \text{ Radial basis function (RBF)}$$

$x \in \mathbb{R}$ Datapoint

$\mu \in \mathbb{R}$ Center (or mean)

$\sigma^2 \in \mathbb{R}$ Variance

Parameters $\{\mu, \sigma^2\}$

Multivariate Gaussian distribution:

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

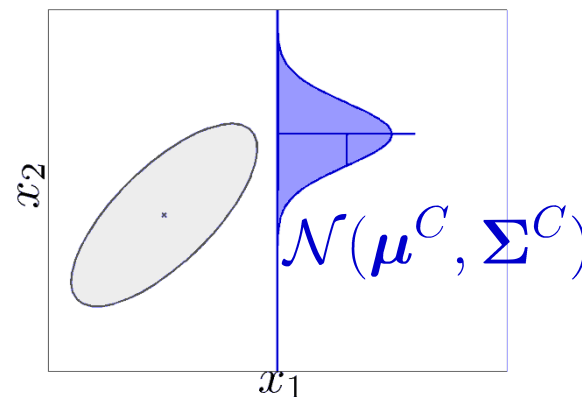
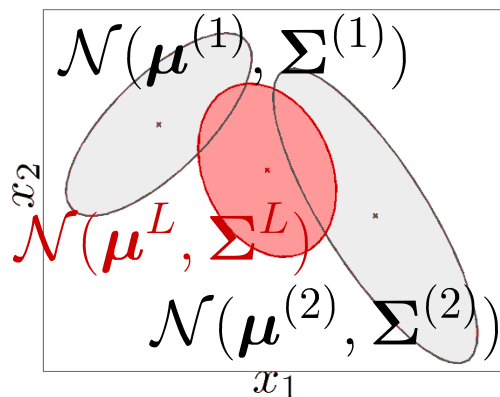
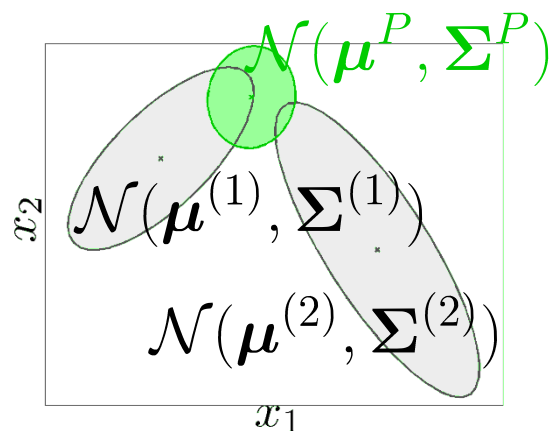
$\mathbf{x} \in \mathbb{R}^D$ Datapoint

$\boldsymbol{\mu} \in \mathbb{R}^D$ Center (or mean)

$\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$ Covariance matrix

Parameters $\{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$

Properties of Gaussian distributions



Linear combination:

$$\mathcal{N}(\boldsymbol{\mu}^L, \boldsymbol{\Sigma}^L) \sim \frac{1}{2} \mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)}) + \frac{1}{2} \mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)})$$

Product of Gaussians:

$$c \mathcal{N}(\boldsymbol{\mu}^P, \boldsymbol{\Sigma}^P) \sim \mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)}) \cdot \mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)})$$

Conditional probability:

$$\mathcal{N}(\boldsymbol{\mu}^C, \boldsymbol{\Sigma}^C) \sim \mathcal{P}(\boldsymbol{x}_2 | \boldsymbol{x}_1)$$

Linear combination

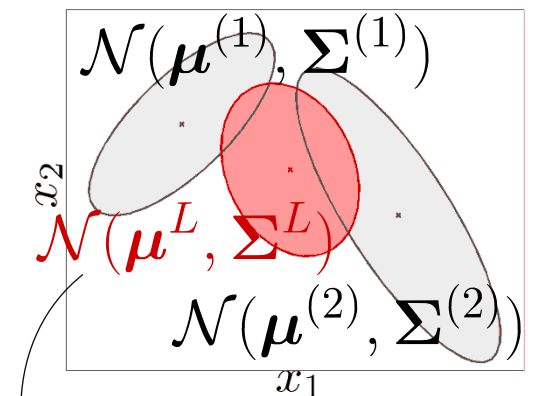
If $\mathbf{x}^{(1)} \sim \mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)})$ and $\mathbf{x}^{(2)} \sim \mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)})$,
the linear transformation $\mathbf{A}^{(1)}\mathbf{x}^{(1)} + \mathbf{A}^{(2)}\mathbf{x}^{(2)} + \mathbf{c}$
follows the distribution

$$\mathbf{A}^{(1)}\mathbf{x}^{(1)} + \mathbf{A}^{(2)}\mathbf{x}^{(2)} + \mathbf{c} \sim \mathcal{N}(\boldsymbol{\mu}^L, \boldsymbol{\Sigma}^L),$$

with

$$\boldsymbol{\mu}^L = \mathbf{A}^{(1)}\boldsymbol{\mu}^{(1)} + \mathbf{A}^{(2)}\boldsymbol{\mu}^{(2)} + \mathbf{c},$$

$$\boldsymbol{\Sigma}^L = \mathbf{A}^{(1)}\boldsymbol{\Sigma}^{(1)}\mathbf{A}^{(1)\top} + \mathbf{A}^{(2)}\boldsymbol{\Sigma}^{(2)}\mathbf{A}^{(2)\top}.$$



for $\mathbf{A}^{(1)} = \mathbf{A}^{(2)} = \frac{1}{2}\mathbf{I}$
and $\mathbf{c} = \mathbf{0}$

Product of Gaussians

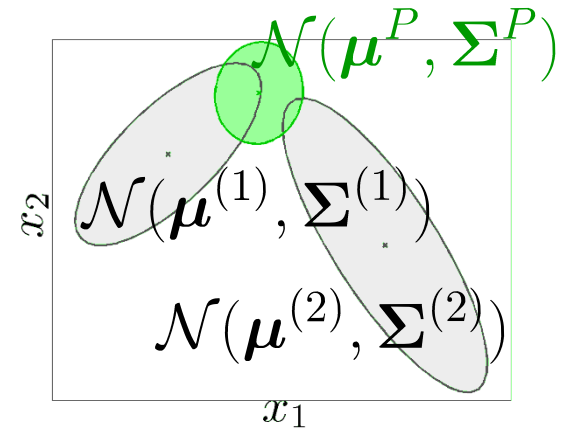
The product of two Gaussian distributions $\mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)})$ and $\mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)})$ is defined by

$$c \mathcal{N}(\boldsymbol{\mu}^P, \boldsymbol{\Sigma}^P) = \mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)}) \cdot \mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)}),$$

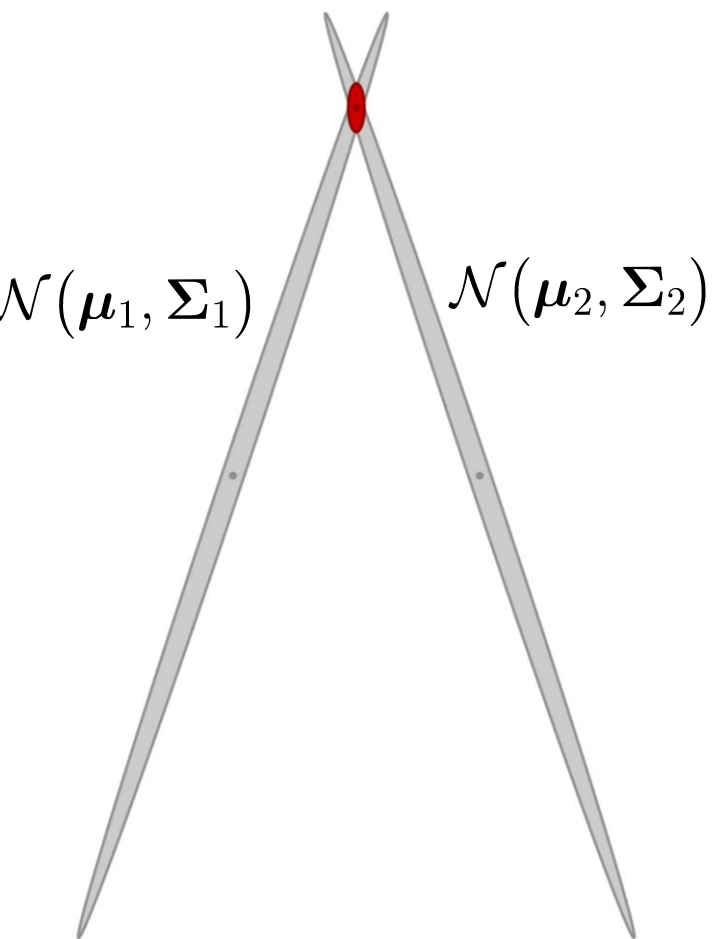
with $c = \mathcal{N}(\boldsymbol{\mu}^{(1)} | \boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(1)} + \boldsymbol{\Sigma}^{(2)}),$

$$\boldsymbol{\Sigma}^P = \left(\boldsymbol{\Sigma}^{(1)-1} + \boldsymbol{\Sigma}^{(2)-1} \right)^{-1},$$

$$\boldsymbol{\mu}^P = \boldsymbol{\Sigma}^P \left(\boldsymbol{\Sigma}^{(1)-1} \boldsymbol{\mu}^{(1)} + \boldsymbol{\Sigma}^{(2)-1} \boldsymbol{\mu}^{(2)} \right).$$



Product of Gaussians



μ_i center of Gaussian

Σ_i covariance matrix

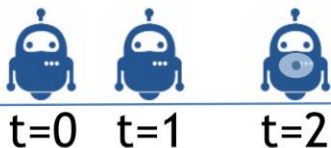
\mathbf{W}_i precision matrix
($\mathbf{W}_i = \Sigma_i^{-1}$)

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mu_1 - \mathbf{x}\|_{\mathbf{W}_1}^2 + \|\mu_2 - \mathbf{x}\|_{\mathbf{W}_2}^2$$

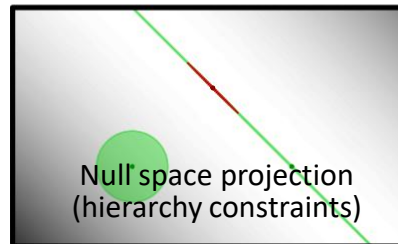
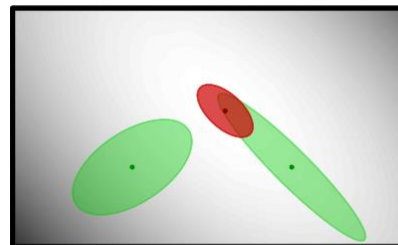
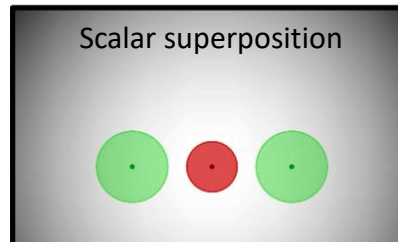
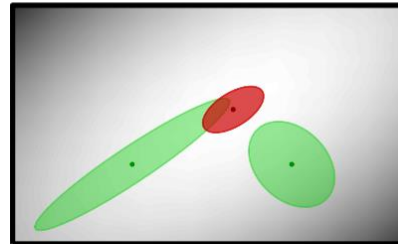
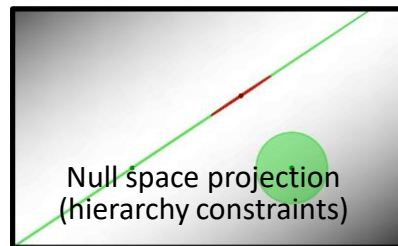
$$= (\mathbf{W}_1 + \mathbf{W}_2)^{-1} (\mathbf{W}_1 \mu_1 + \mathbf{W}_2 \mu_2)$$

→ **Product of Gaussians**

Standard fusion problem
for state estimation



PoG can cover a wide range of behaviors!



Conditional probability

Let $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be defined by

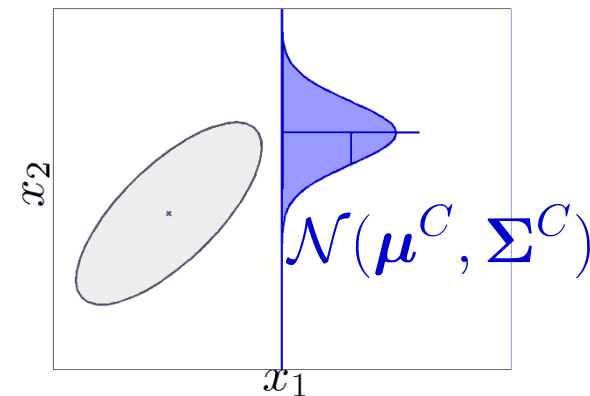
$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

The conditional probability $\mathcal{P}(\mathbf{x}_2|\mathbf{x}_1)$ is defined by

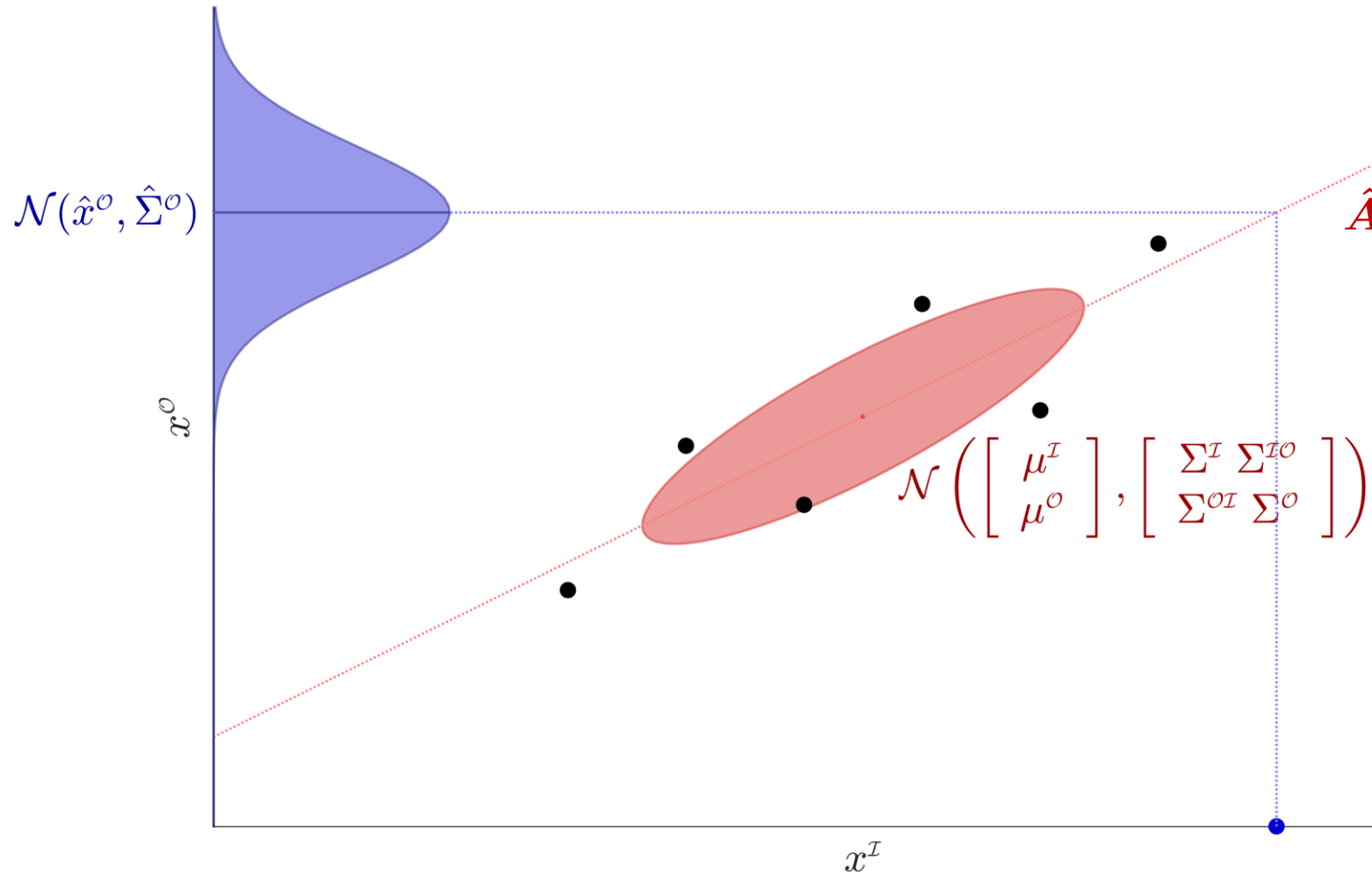
$$\mathcal{P}(\mathbf{x}_2|\mathbf{x}_1) \sim \mathcal{N}(\boldsymbol{\mu}^C, \boldsymbol{\Sigma}^C),$$

with

$$\begin{aligned} \boldsymbol{\mu}^C &= \mu_2 + \boldsymbol{\Sigma}_{21}(\boldsymbol{\Sigma}_{11})^{-1}(\mathbf{x}_1 - \mu_1), \\ \boldsymbol{\Sigma}^C &= \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}(\boldsymbol{\Sigma}_{11})^{-1}\boldsymbol{\Sigma}_{12}. \end{aligned}$$



Conditional distribution



$$\begin{aligned} \hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} (\mathbf{X}^O - \mathbf{X}^I \mathbf{A})^\top (\mathbf{X}^O - \mathbf{X}^I \mathbf{A}) \\ &= (\mathbf{X}^{I\top} \mathbf{X}^I)^{-1} \mathbf{X}^{I\top} \mathbf{X}^O = \mathbf{X}^{I\dagger} \mathbf{X}^O \end{aligned}$$

→ Linear regression from joint distribution

Conditional distribution

We consider multivariate datapoints \mathbf{x} and multivariate Gaussian distributions characterized by centers $\boldsymbol{\mu}$ and covariances $\boldsymbol{\Sigma}$, that can be partitioned as

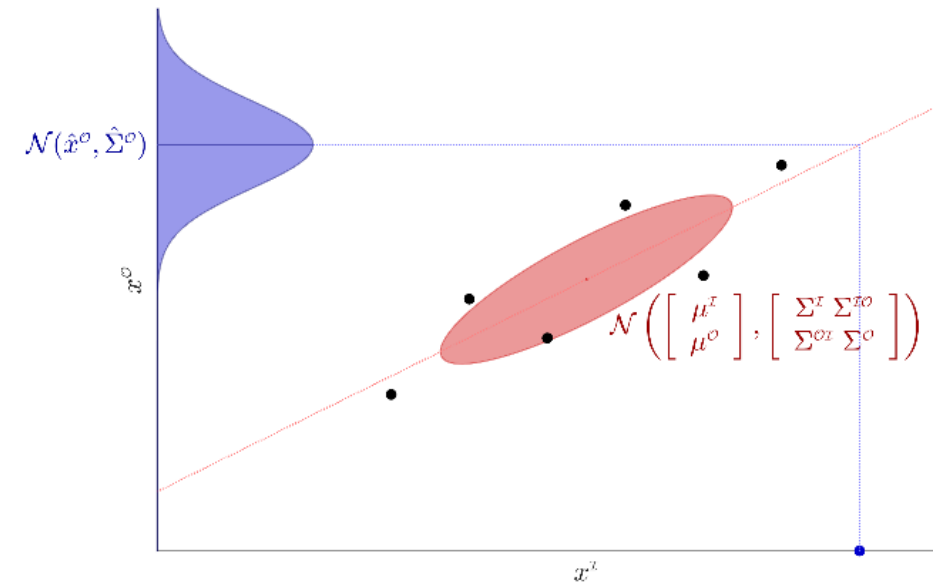
$$\mathbf{x} = \begin{bmatrix} \mathbf{x}^I \\ \mathbf{x}^O \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}^I \\ \boldsymbol{\mu}^O \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}^I & \boldsymbol{\Sigma}^{IO} \\ \boldsymbol{\Sigma}^{OI} & \boldsymbol{\Sigma}^O \end{bmatrix}.$$

If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we have $\mathbf{x}^O | \mathbf{x}^I \sim \mathcal{N}(\hat{\mathbf{x}}^O, \hat{\boldsymbol{\Sigma}}^O)$, with parameters

$$\begin{aligned} \hat{\mathbf{x}}^O &= \boldsymbol{\mu}^O + \boldsymbol{\Sigma}^{OI} \boldsymbol{\Sigma}^{I-1} (\mathbf{x}^I - \boldsymbol{\mu}^I), \\ \hat{\boldsymbol{\Sigma}}^O &= \boldsymbol{\Sigma}^O - \boldsymbol{\Sigma}^{OI} \boldsymbol{\Sigma}^{I-1} \boldsymbol{\Sigma}^{IO}. \end{aligned}$$

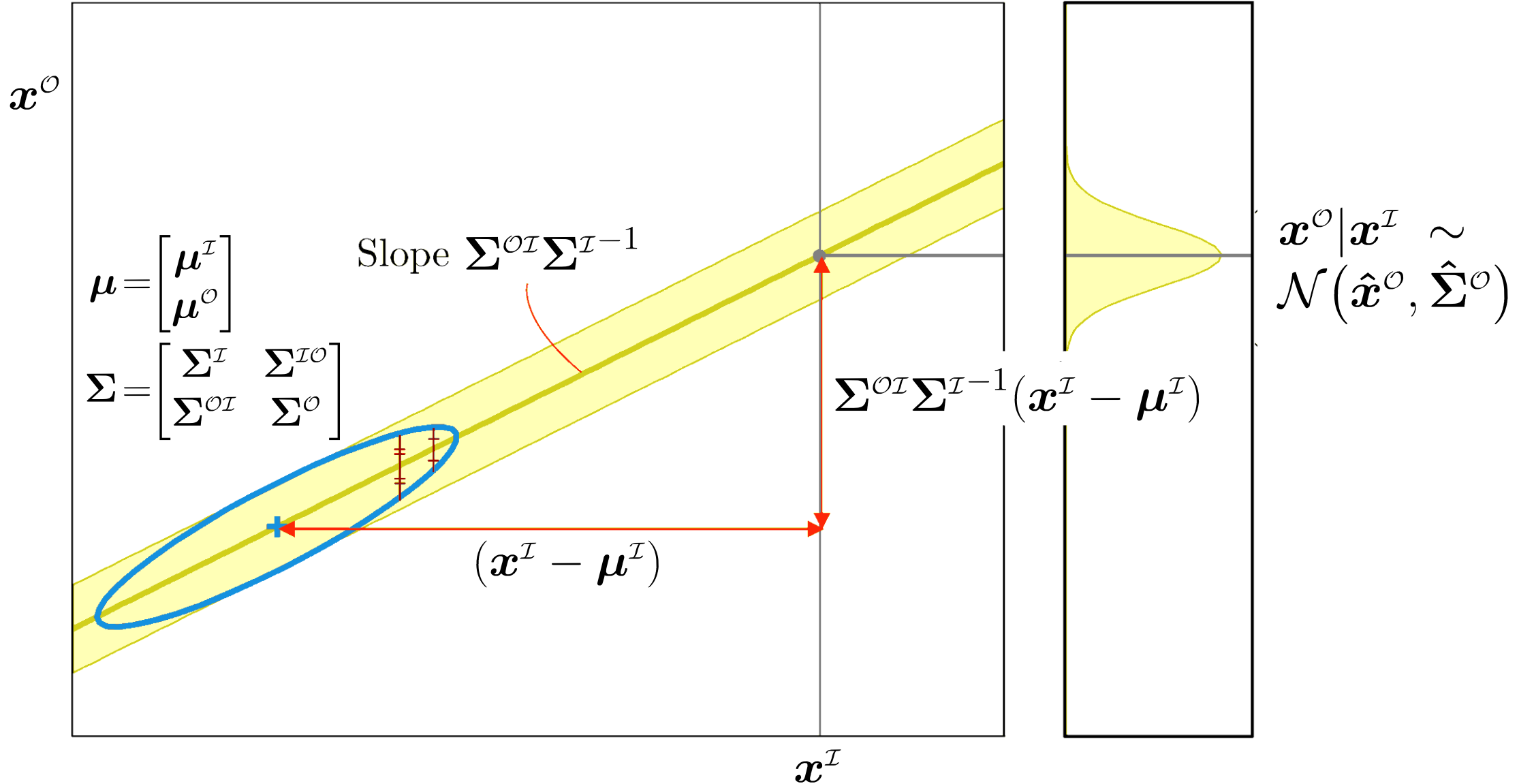
We can see that $\hat{\mathbf{x}}^O$ is linearly dependent on \mathbf{x}^I , and that $\hat{\boldsymbol{\Sigma}}^O$ is independent of \mathbf{x}^I .

We can also notice that for full joint covariance, the conditional covariance $\hat{\boldsymbol{\Sigma}}^O$ will typically be smaller than the marginal $\boldsymbol{\Sigma}^O$.



Conditional distribution

$$\hat{x}^o = \mu^o + \Sigma^{oI} \Sigma^{I-1} (x^I - \mu^I)$$



Locally weighted regression (LWR)

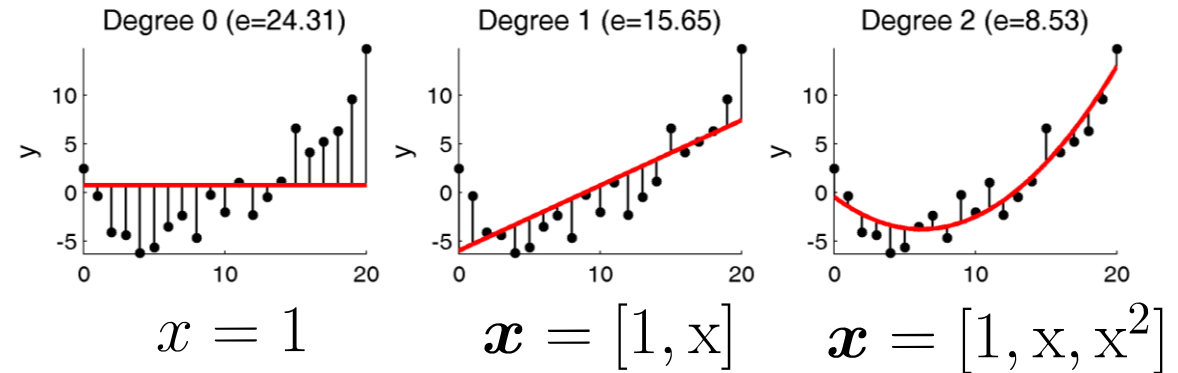
Python notebooks:
demo_LWR.ipynb

Matlab codes:
demo_LWR01.m

Recap: Linear regression

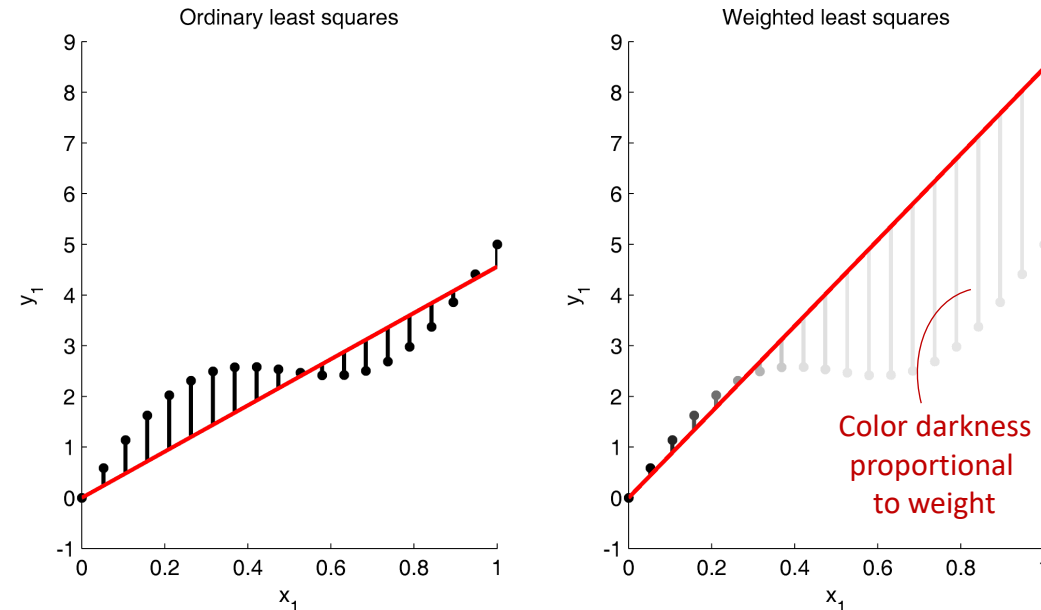
$$\hat{A} = \arg \min_A (Y - XA)^T (Y - XA)$$

$$= (X^T X)^{-1} X^T Y = X^\dagger Y$$



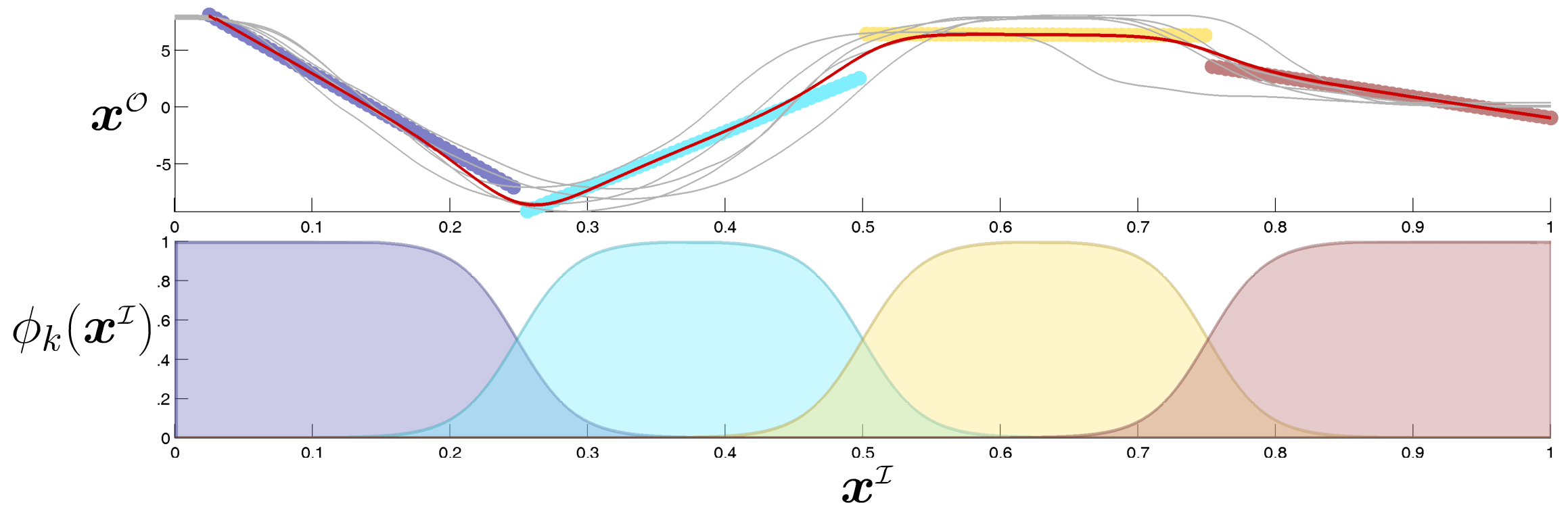
$$\hat{A} = \arg \min_A (Y - XA)^T W (Y - XA)$$

$$= (X^T W X)^{-1} X^T W Y$$



Locally weighted regression (LWR)

K weighted regressions with different \mathbf{W} are performed on the same dataset $\{\mathbf{X}^I, \mathbf{X}^O\}$



Locally weighted regression (LWR)

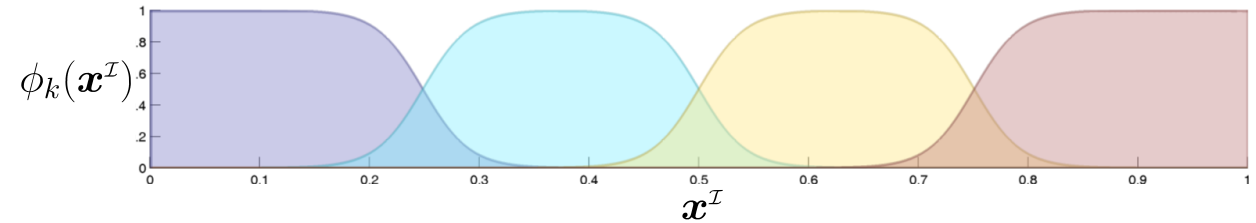
LWR computes K estimates $\hat{\mathbf{A}}_k$, by using different weighting function $\phi_k(\mathbf{x}_n^{\mathcal{I}})$, often defined as the **radial basis functions** (RBF)

$$\tilde{\phi}_k(\mathbf{x}_n^{\mathcal{I}}) = \exp\left(-\frac{1}{2}(\mathbf{x}_n^{\mathcal{I}} - \boldsymbol{\mu}_k^{\mathcal{I}})^{\top} \boldsymbol{\Sigma}_k^{\mathcal{I}}^{-1}(\mathbf{x}_n^{\mathcal{I}} - \boldsymbol{\mu}_k^{\mathcal{I}})\right),$$

or in a rescaled form as

$$\phi_k(\mathbf{x}_n^{\mathcal{I}}) = \frac{\tilde{\phi}_k(\mathbf{x}_n^{\mathcal{I}})}{\sum_{i=1}^K \tilde{\phi}_i(\mathbf{x}_n^{\mathcal{I}})},$$

where $\boldsymbol{\mu}_k^{\mathcal{I}}$ and $\boldsymbol{\Sigma}_k^{\mathcal{I}}$ are the parameters of the k -th RBF.



→ **Nonlinear problem solved locally by linear regression**

Locally weighted regression (LWR)

Often, the centroids $\boldsymbol{\mu}_k^{\mathcal{I}}$ are set to uniformly cover the input space, and $\boldsymbol{\Sigma}_k^{\mathcal{I}} = \mathbf{I}\sigma^2$ is used as a common bandwidth shared by all basis functions.

An associated diagonal matrix

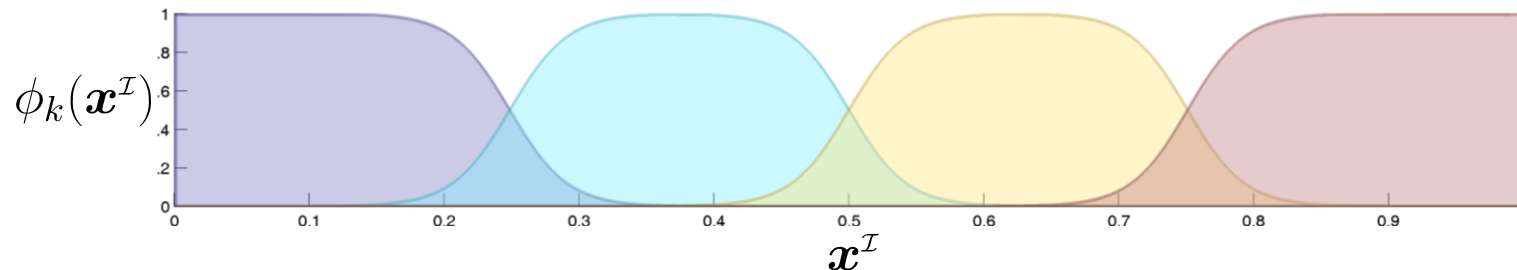
$$\mathbf{W}_k = \text{diag}\left(\phi_k(\mathbf{x}_1^{\mathcal{I}}), \phi_k(\mathbf{x}_2^{\mathcal{I}}), \dots, \phi_k(\mathbf{x}_N^{\mathcal{I}})\right)$$

can be used to evaluate $\hat{\mathbf{A}}_k$. The result can then be used to compute

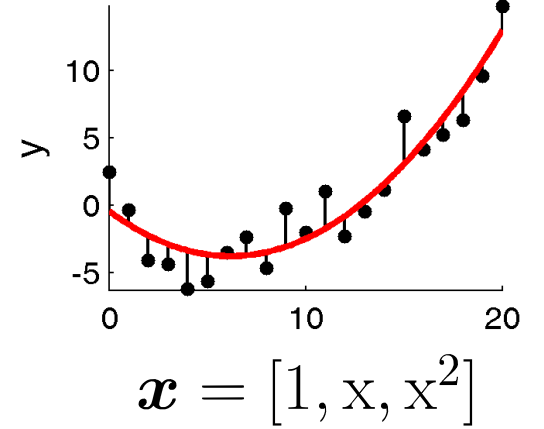
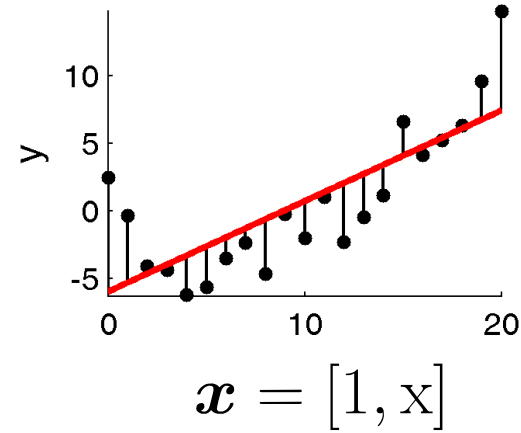
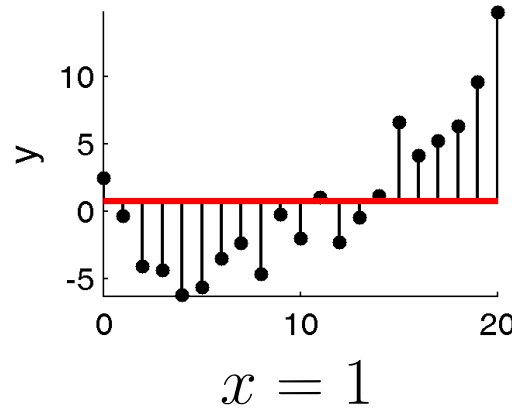
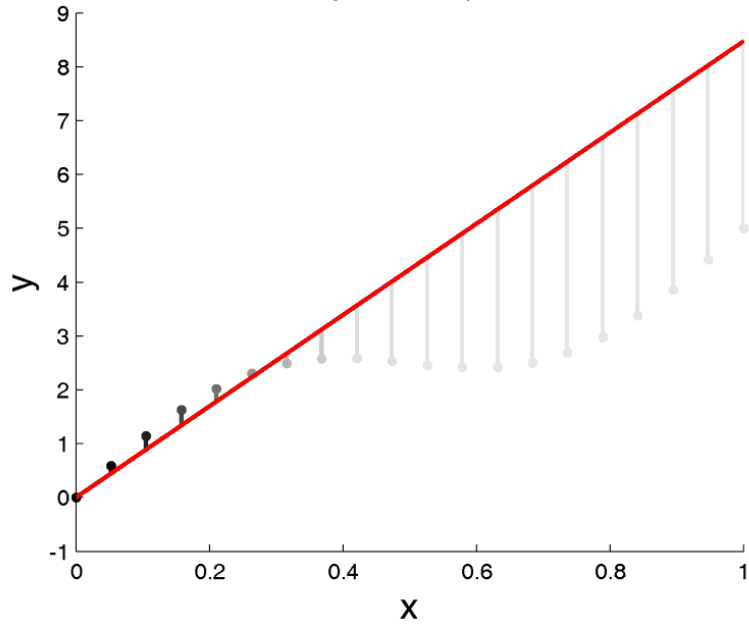
$$\mathbf{X}^{\mathcal{I}} = [t_1, t_2, \dots, t_N]^{\top}$$

$$\hat{\mathbf{A}}_k = (\mathbf{X}^{\mathcal{I}\top} \mathbf{W}_k \mathbf{X}^{\mathcal{I}})^{-1} \mathbf{X}^{\mathcal{I}\top} \mathbf{W}_k \mathbf{X}^{\mathcal{O}}$$

$$\mathbf{X}^{\mathcal{O}} = \sum_{k=1}^K \mathbf{W}_k \mathbf{X}^{\mathcal{I}} \hat{\mathbf{A}}_k$$



Locally weighted regression (LWR)

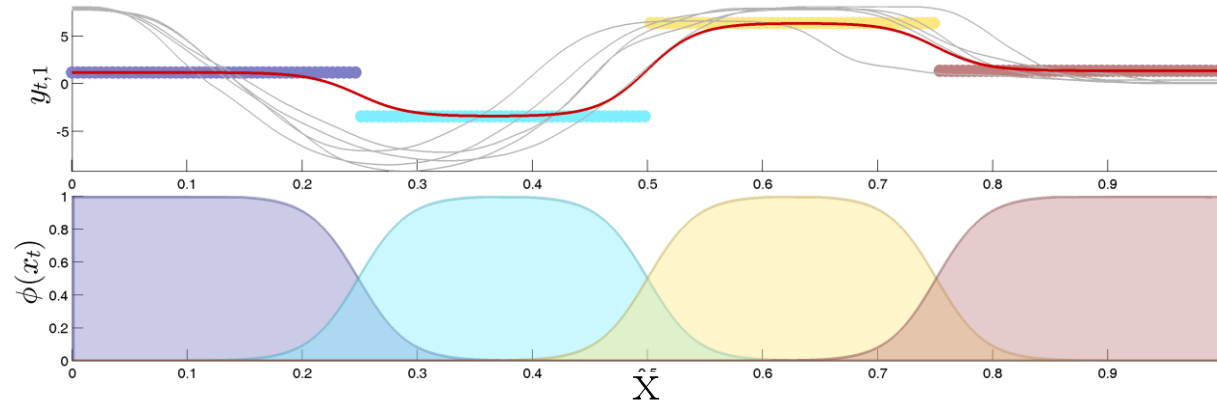
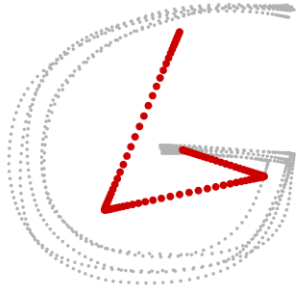


$$\hat{A} = (X^T W X)^{-1} X^T W Y$$

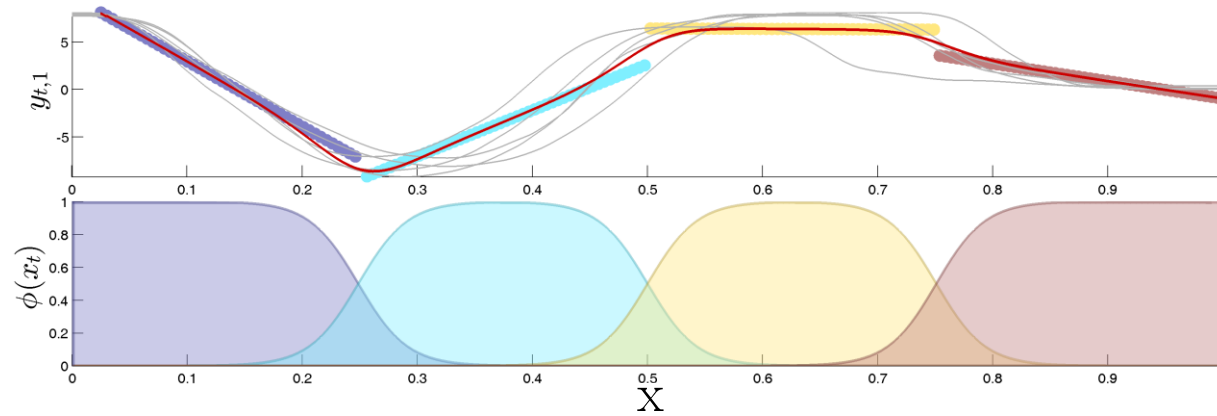
LWR can be used for local least squares polynomial fitting by changing the definition of the inputs.

Locally weighted regression (LWR)

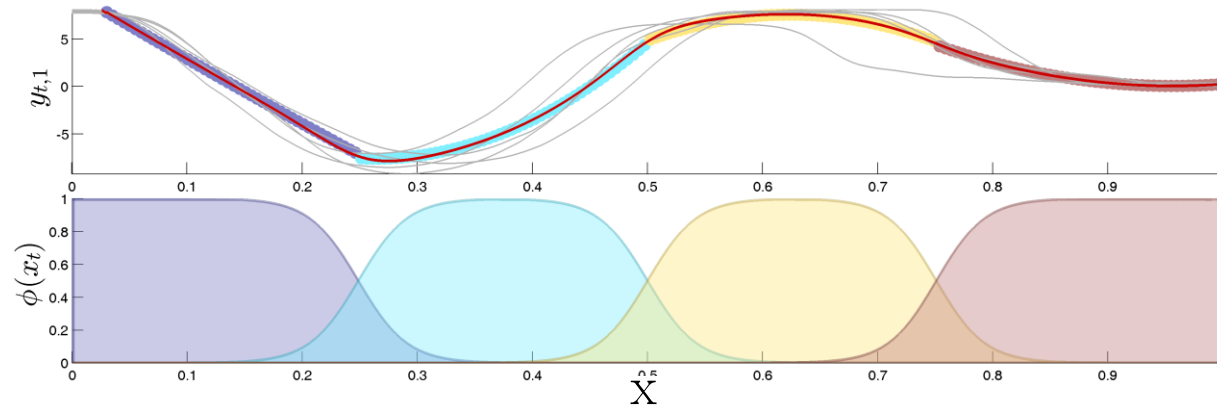
$$\mathbf{x} = 1$$



$$\mathbf{x} = [1, x]$$

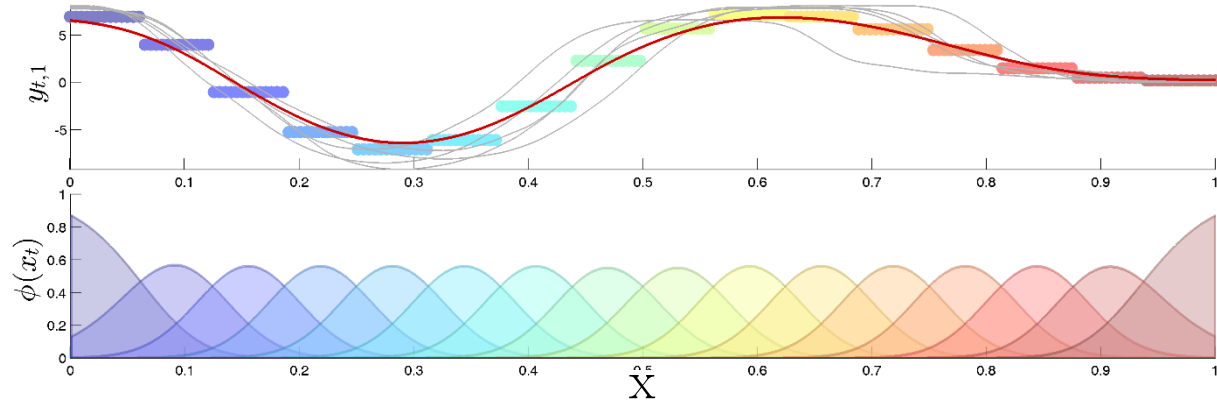


$$\mathbf{x} = [1, x, x^2]$$

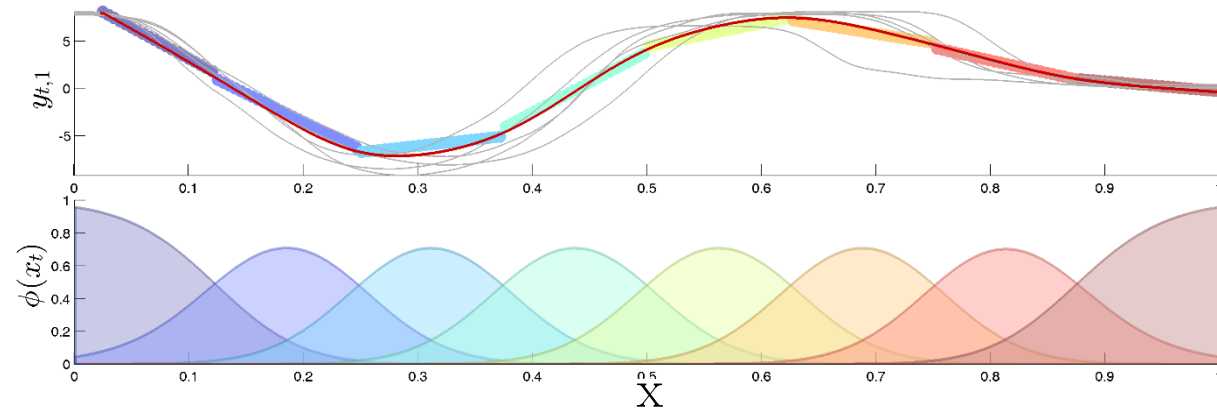


Locally weighted regression (LWR)

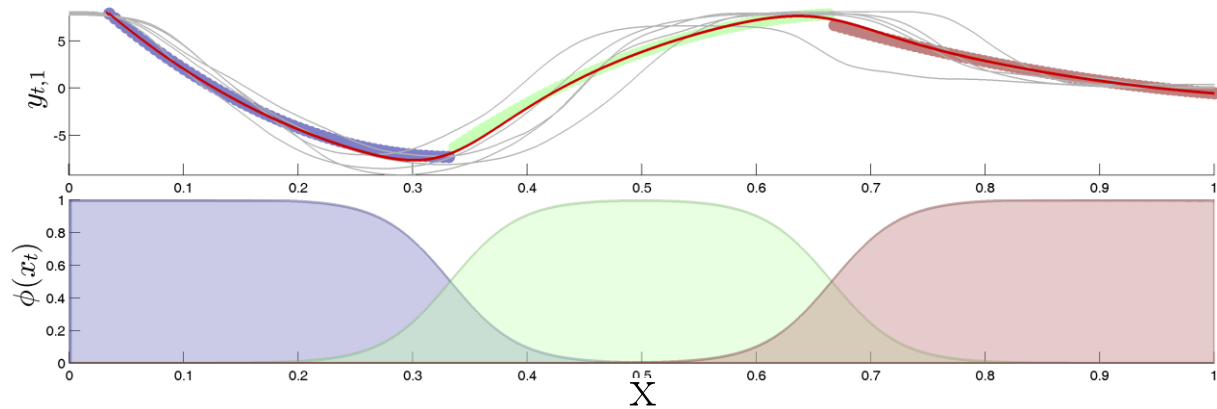
$$\mathbf{x} = 1$$



$$\mathbf{x} = [1, x]$$



$$\mathbf{x} = [1, x, x^2]$$



Gaussian mixture regression (GMR)

Python notebooks:
demo_GMR.ipynb

Matlab codes:
demo_GMR01.m
demo_GMR_polyFit01.m

Gaussian Mixture Model (GMM)

$$\mathbf{x}_n \sim \sum_{i=1}^K \pi_i \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

$$\mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_i)\right)$$

K Gaussians
 N datapoints of dimension D

$\mathbf{x}_n \in \mathbb{R}^D$ Datapoint

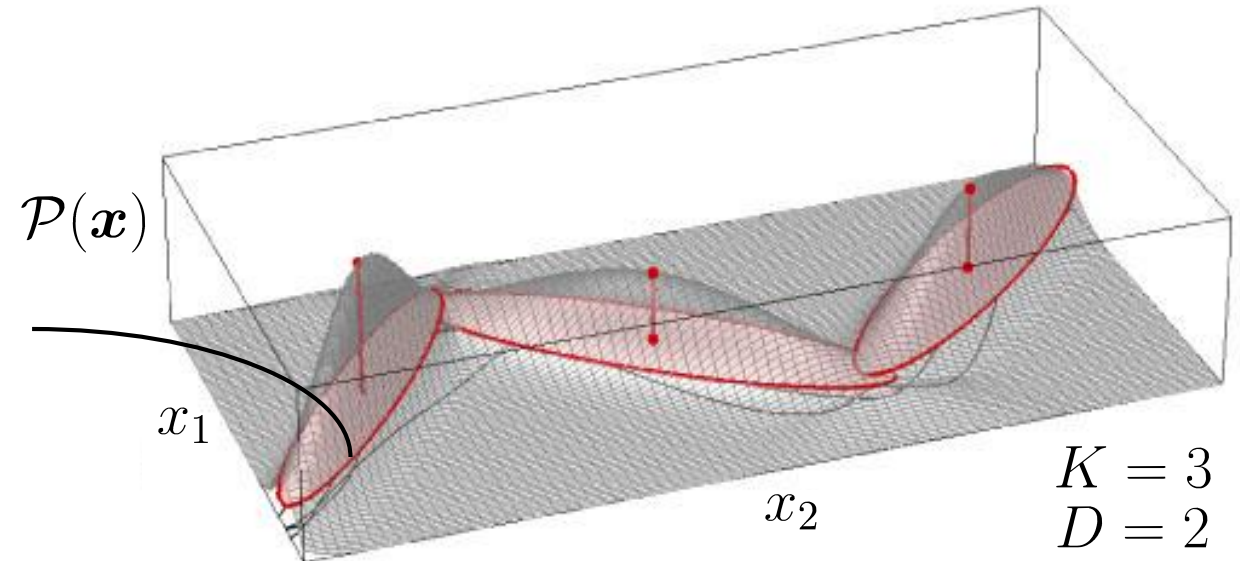
$\pi_i \in \mathbb{R}$ Mixing coefficient

$\boldsymbol{\mu}_i \in \mathbb{R}^D$ Center (or mean)

$\boldsymbol{\Sigma}_i \in \mathbb{R}^{D \times D}$ Covariance matrix

Parameters $\Theta^{\text{GMM}} = \{\pi_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^K$

Equidensity contour of
one standard deviation



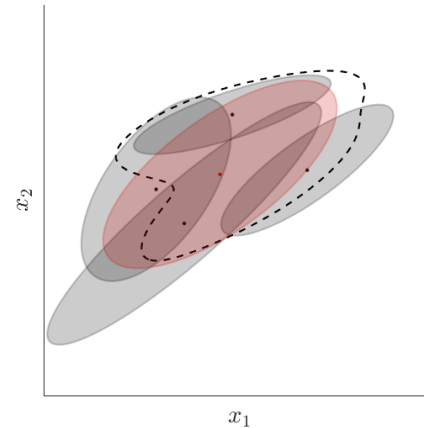
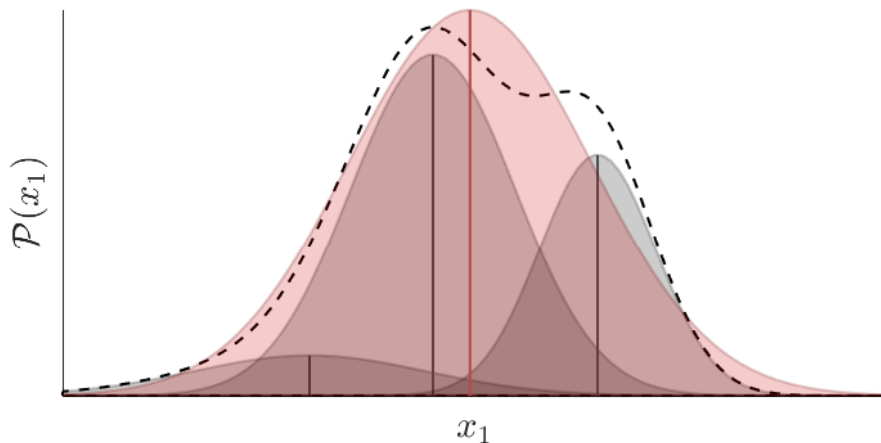
Gaussian estimate of a mixture of Gaussians

We can approximate a mixture of Gaussians $\sum_{i=1}^K h_i \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ with a single Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, by **moment matching of the means (first moments) and covariances (second moments)** with

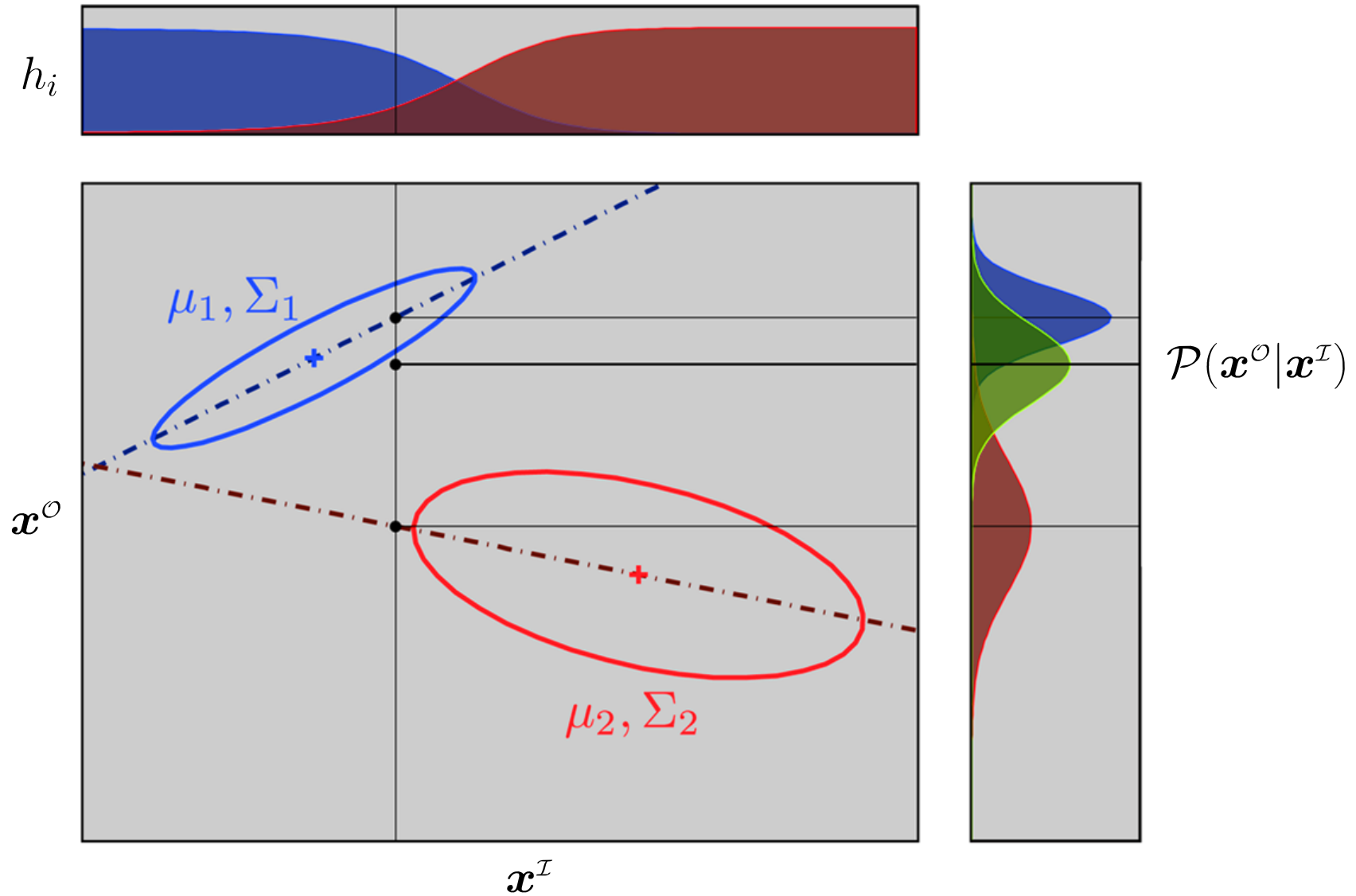
$$\boldsymbol{\mu} = \sum_{i=1}^K h_i \boldsymbol{\mu}_i,$$

$$\boldsymbol{\Sigma} = \sum_{i=1}^K h_i \left(\boldsymbol{\Sigma}_i + \boldsymbol{\mu}_i \boldsymbol{\mu}_i^\top \right) - \boldsymbol{\mu} \boldsymbol{\mu}^\top,$$

also referred to as the **law of total mean and (co)variance**.



Gaussian mixture regression (GMR)



Gaussian mixture regression (GMR)

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}^I \\ \mathbf{x}^O \end{bmatrix} \quad \boldsymbol{\mu}_i = \begin{bmatrix} \boldsymbol{\mu}_i^I \\ \boldsymbol{\mu}_i^O \end{bmatrix} \quad \boldsymbol{\Sigma}_i = \begin{bmatrix} \boldsymbol{\Sigma}_i^I & \boldsymbol{\Sigma}_i^{IO} \\ \boldsymbol{\Sigma}_i^{OI} & \boldsymbol{\Sigma}_i^O \end{bmatrix}$$

$\mathcal{P}(\mathbf{x}^O | \mathbf{x}^I)$ can be computed as the multimodal conditional distribution

$$\mathcal{P}(\mathbf{x}^O | \mathbf{x}^I) = \sum_{i=1}^K h_i \mathcal{N}(\mathbf{x}^O | \hat{\boldsymbol{\mu}}_i^O, \hat{\boldsymbol{\Sigma}}_i^O),$$

$$\text{with } \hat{\boldsymbol{\mu}}_i^O = \boldsymbol{\mu}_i^O + \boldsymbol{\Sigma}_i^{OI} \boldsymbol{\Sigma}_i^I{}^{-1} (\mathbf{x}^I - \boldsymbol{\mu}_i^I),$$

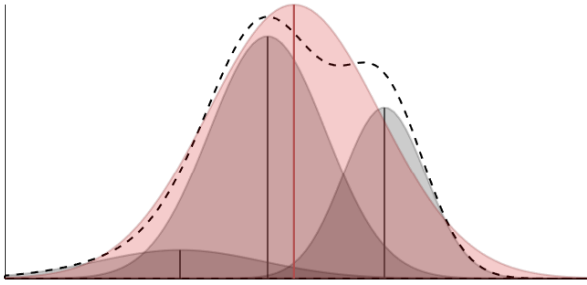
$$\hat{\boldsymbol{\Sigma}}_i^O = \boldsymbol{\Sigma}_i^O - \boldsymbol{\Sigma}_i^{OI} \boldsymbol{\Sigma}_i^I{}^{-1} \boldsymbol{\Sigma}_i^{IO}$$

$$\text{and } h_i = \frac{\pi_i \mathcal{N}(\mathbf{x}^I | \boldsymbol{\mu}_i^I, \boldsymbol{\Sigma}_i^I)}{\sum_k^K \pi_k \mathcal{N}(\mathbf{x}^I | \boldsymbol{\mu}_k^I, \boldsymbol{\Sigma}_k^I)},$$

computed with the marginal

$$\mathcal{N}(\mathbf{x}^I | \boldsymbol{\mu}_i^I, \boldsymbol{\Sigma}_i^I) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}_i^I|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}^I - \boldsymbol{\mu}_i^I)^\top \boldsymbol{\Sigma}_i^I{}^{-1} (\mathbf{x}^I - \boldsymbol{\mu}_i^I)\right).$$

Gaussian mixture regression (GMR)



$$\mathbf{x} = \begin{bmatrix} \mathbf{x}^I \\ \mathbf{x}^O \end{bmatrix} \quad \boldsymbol{\mu}_i = \begin{bmatrix} \boldsymbol{\mu}_i^I \\ \boldsymbol{\mu}_i^O \end{bmatrix} \quad \boldsymbol{\Sigma}_i = \begin{bmatrix} \boldsymbol{\Sigma}_i^I & \boldsymbol{\Sigma}_i^{IO} \\ \boldsymbol{\Sigma}_i^{OI} & \boldsymbol{\Sigma}_i^O \end{bmatrix}$$

$$\hat{\boldsymbol{\mu}}_i^O = \boldsymbol{\mu}_i^O + \boldsymbol{\Sigma}_i^{OI} \boldsymbol{\Sigma}_i^{I-1} (\mathbf{x}^I - \boldsymbol{\mu}_i^I)$$

$$\hat{\boldsymbol{\Sigma}}_i^O = \boldsymbol{\Sigma}_i^O - \boldsymbol{\Sigma}_i^{OI} \boldsymbol{\Sigma}_i^{I-1} \boldsymbol{\Sigma}_i^{IO}$$

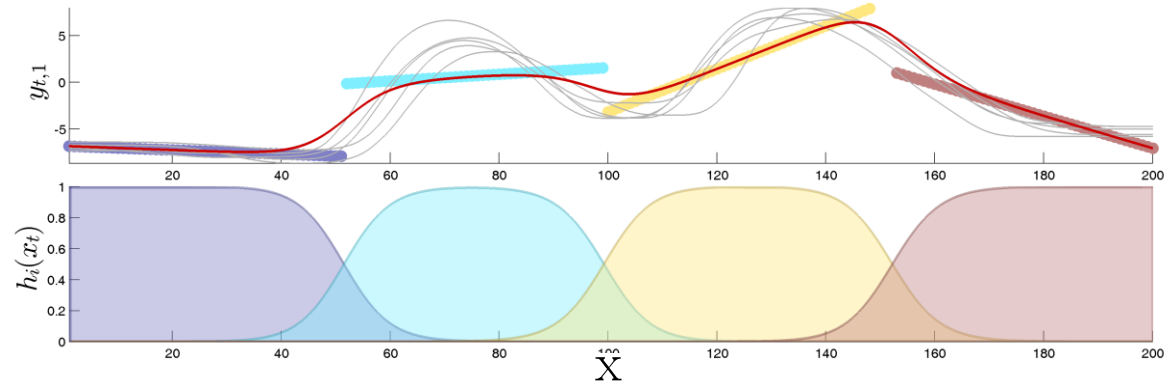
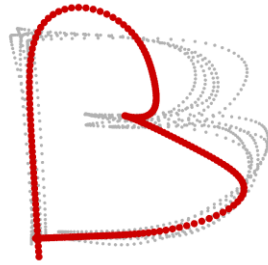
In GMR, an output distribution as a single multivariate Gaussian can be evaluated by moment matching of the means and covariances. The resulting Gaussian distribution $\mathcal{N}(\hat{\boldsymbol{\mu}}^O, \hat{\boldsymbol{\Sigma}}^O)$ has parameters

$$\hat{\boldsymbol{\mu}}^O = \sum_{i=1}^K h_i \hat{\boldsymbol{\mu}}_i^O,$$

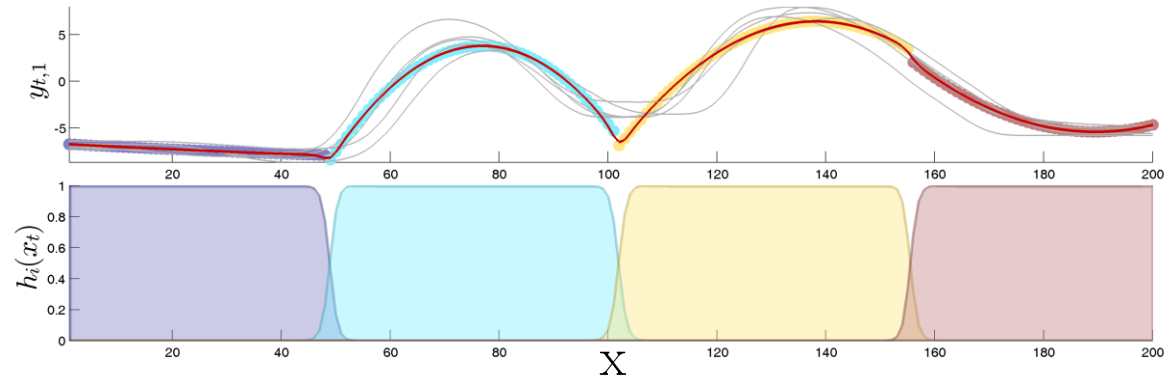
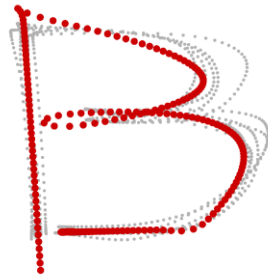
$$\hat{\boldsymbol{\Sigma}}^O = \sum_{i=1}^K h_i \left(\hat{\boldsymbol{\Sigma}}_i^O + \hat{\boldsymbol{\mu}}_i^O \hat{\boldsymbol{\mu}}_i^{O\top} \right) - \hat{\boldsymbol{\mu}}^O \hat{\boldsymbol{\mu}}^{O\top}.$$

GMR for piecewise polynomial fitting

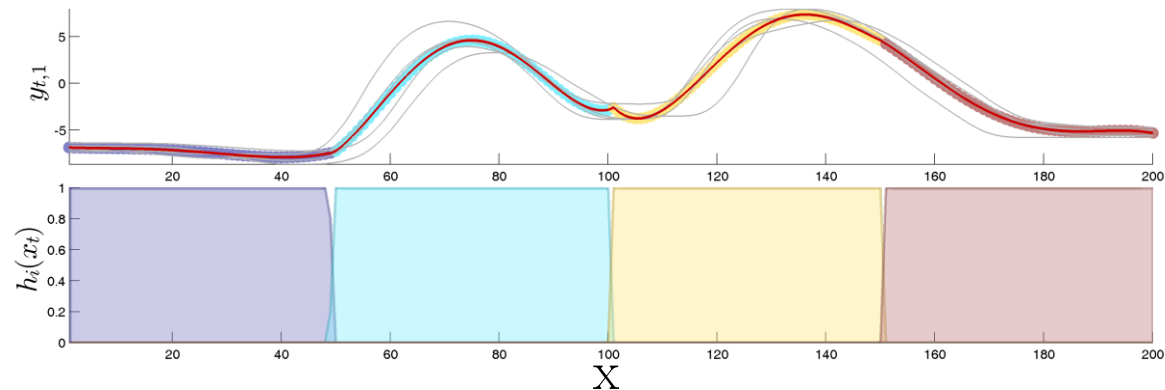
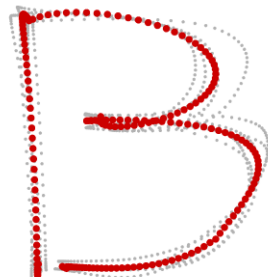
$$\mathbf{x} = 1$$



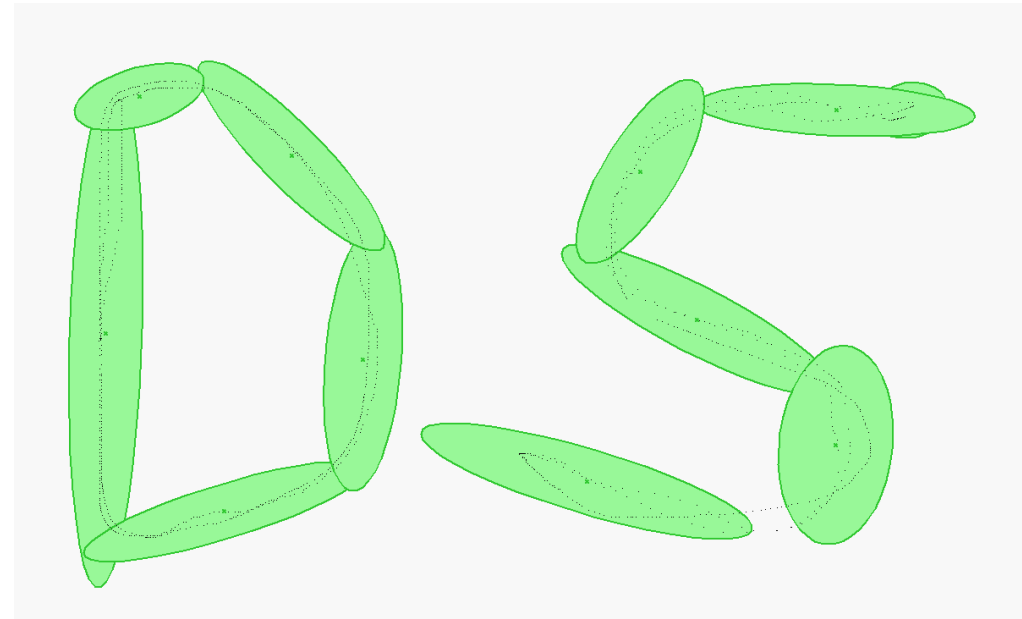
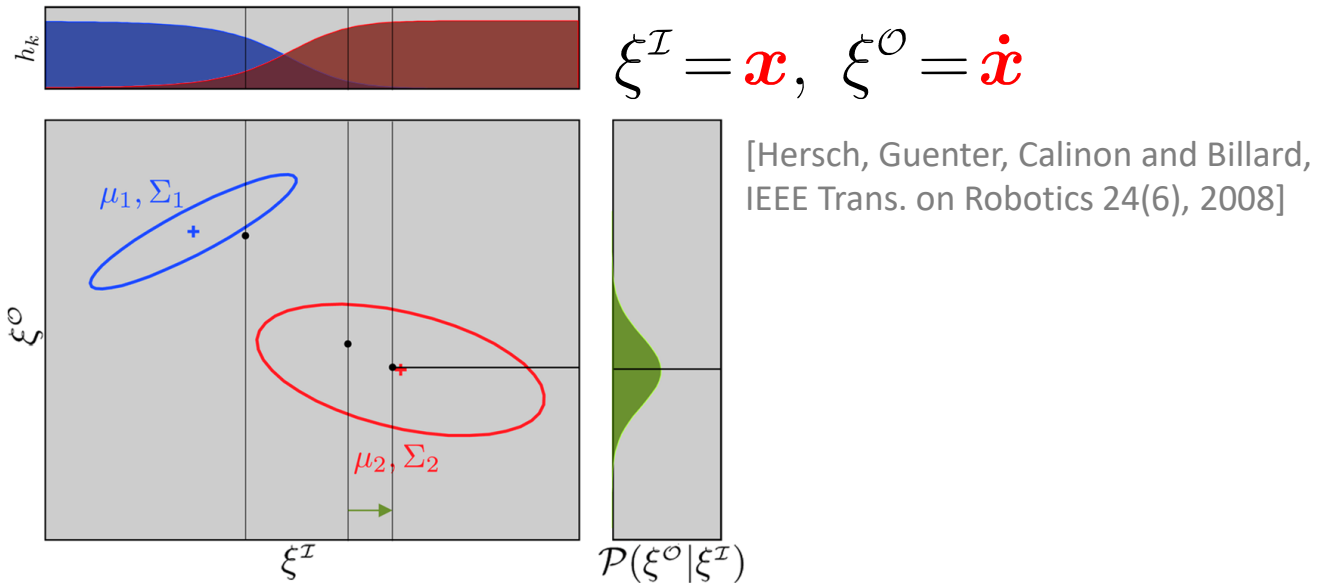
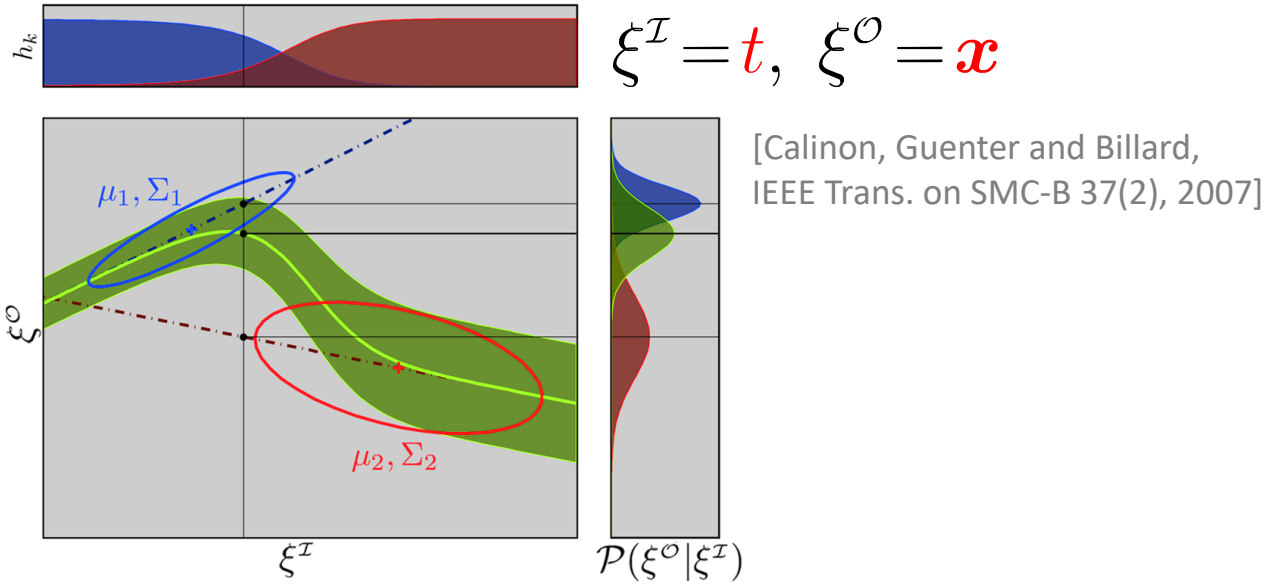
$$\mathbf{x} = [1, \mathbf{x}]$$



$$\mathbf{x} = [1, \mathbf{x}, \mathbf{x}^2]$$



Gaussian mixture regression (GMR)



Gaussian process regression (GPR)

Python notebook:
demo_GPR.ipynb

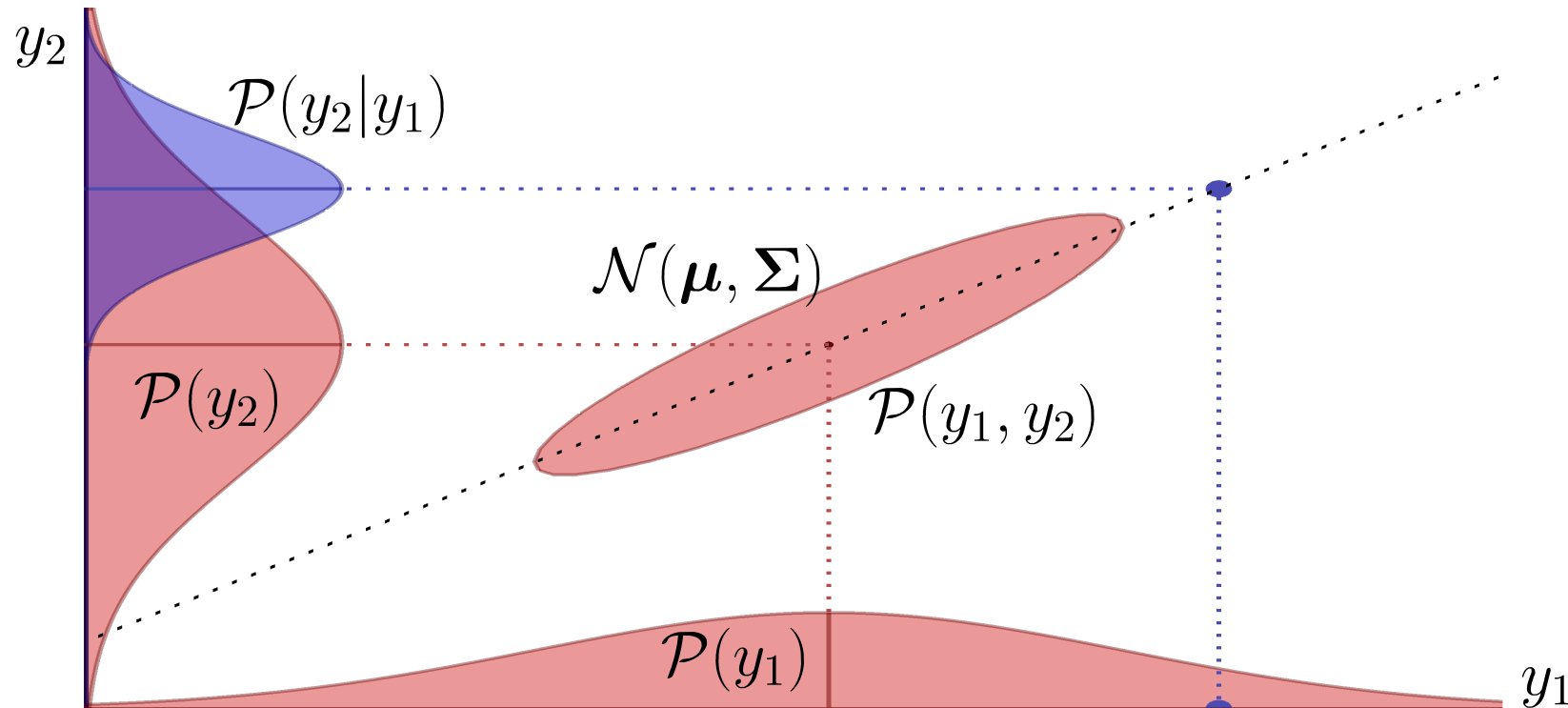
Matlab code:
demo_GPR01.m

Gaussian process - Interpretation

- A joint distribution represented by a bivariate Gaussian forms marginal distributions $P(y_1)$ and $P(y_2)$ that are unidimensional.
- Observing y_1 changes our belief about y_2 , giving rise to a **conditional distribution**.
- Knowledge of the covariance lets us shrink uncertainty in one variable based on the observation of the other.

$$\hat{y}_2 = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(y_1 - \mu_1)$$

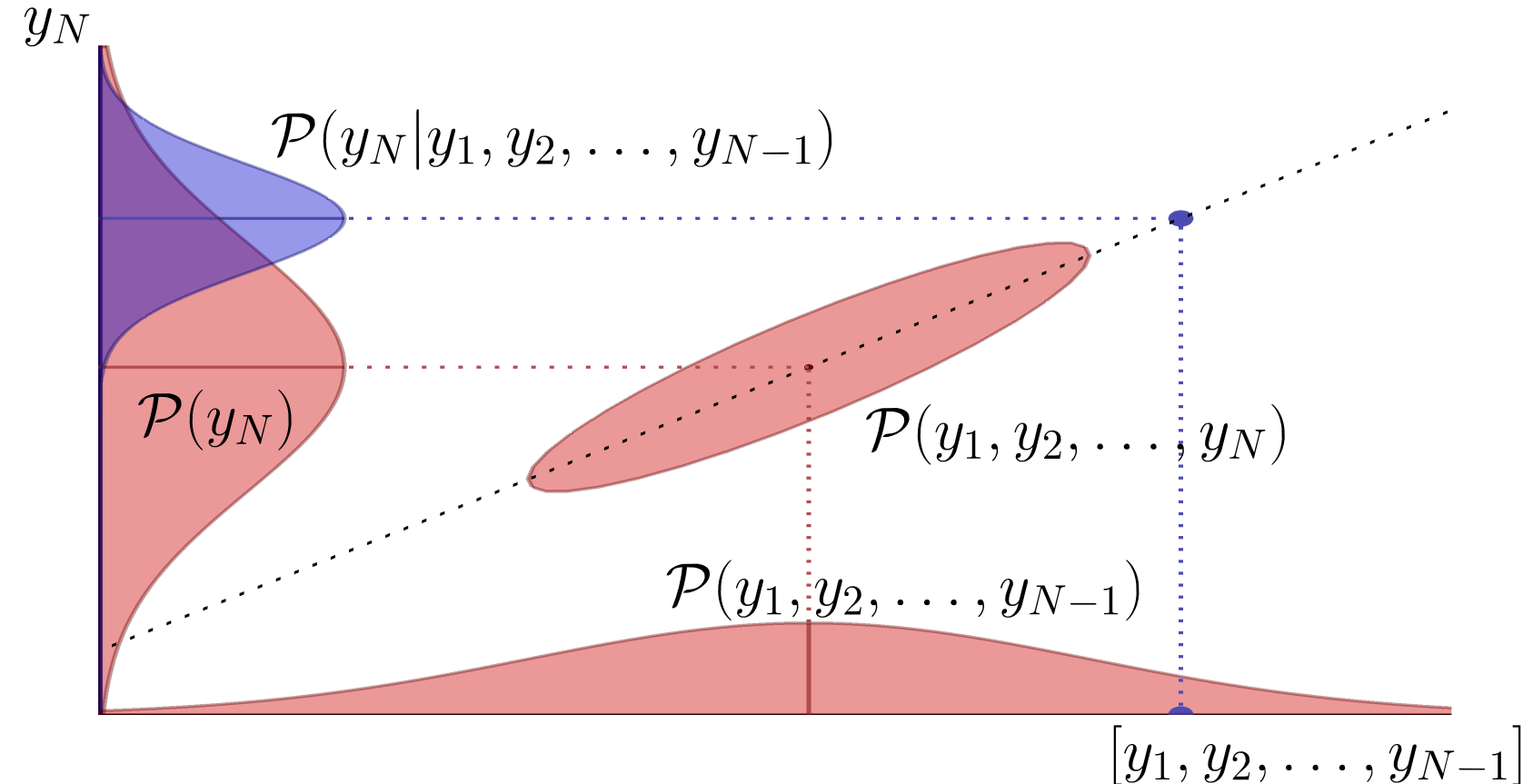
$$\hat{\Sigma}_{22} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$$



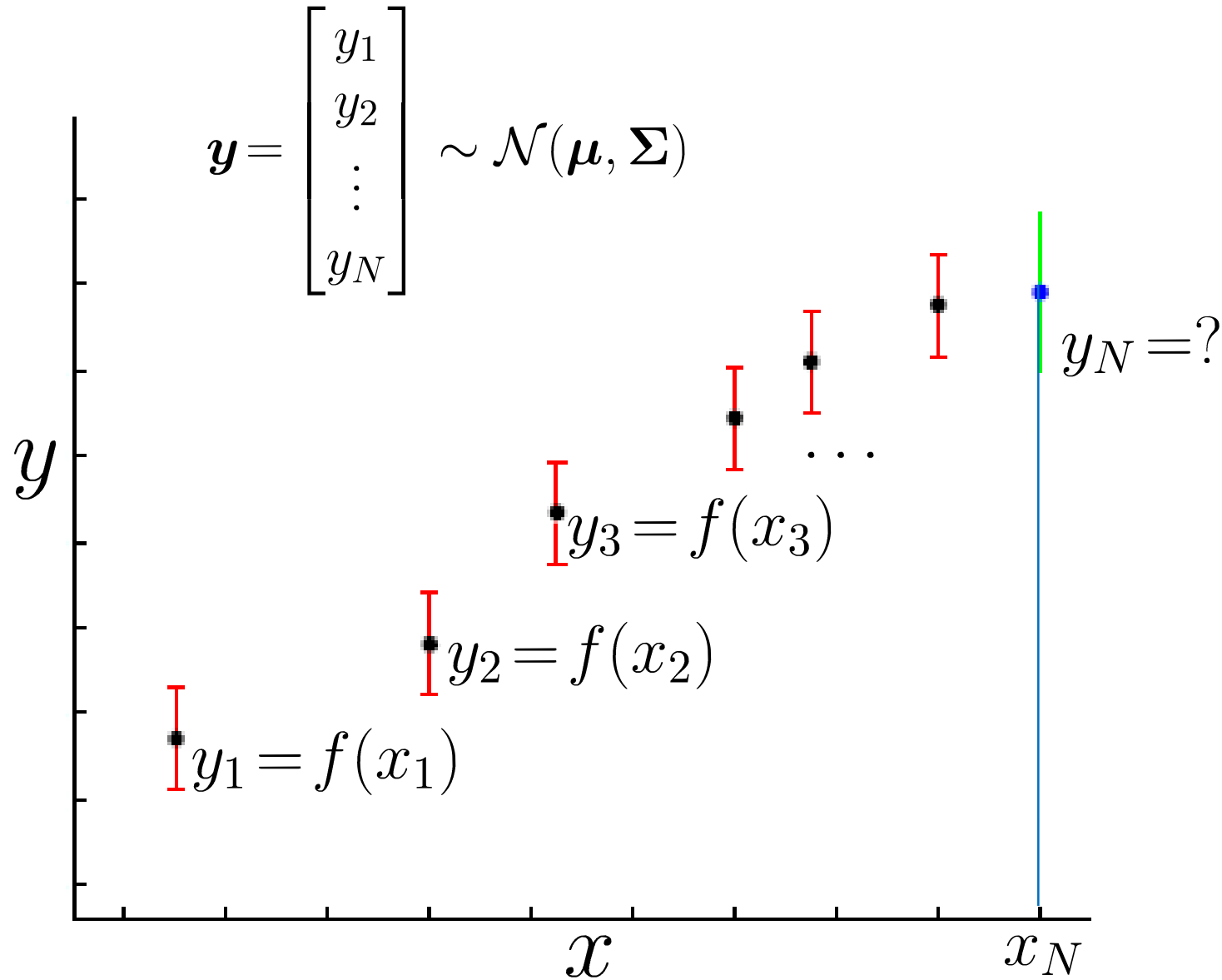
Gaussian process - Interpretation

- This bivariate example can be extended to an arbitrarily large number of variables.
- Indeed, observations in an arbitrary dataset can always be imagined as a single point sampled from a multivariate Gaussian distribution.

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$



Gaussian process - Interpretation

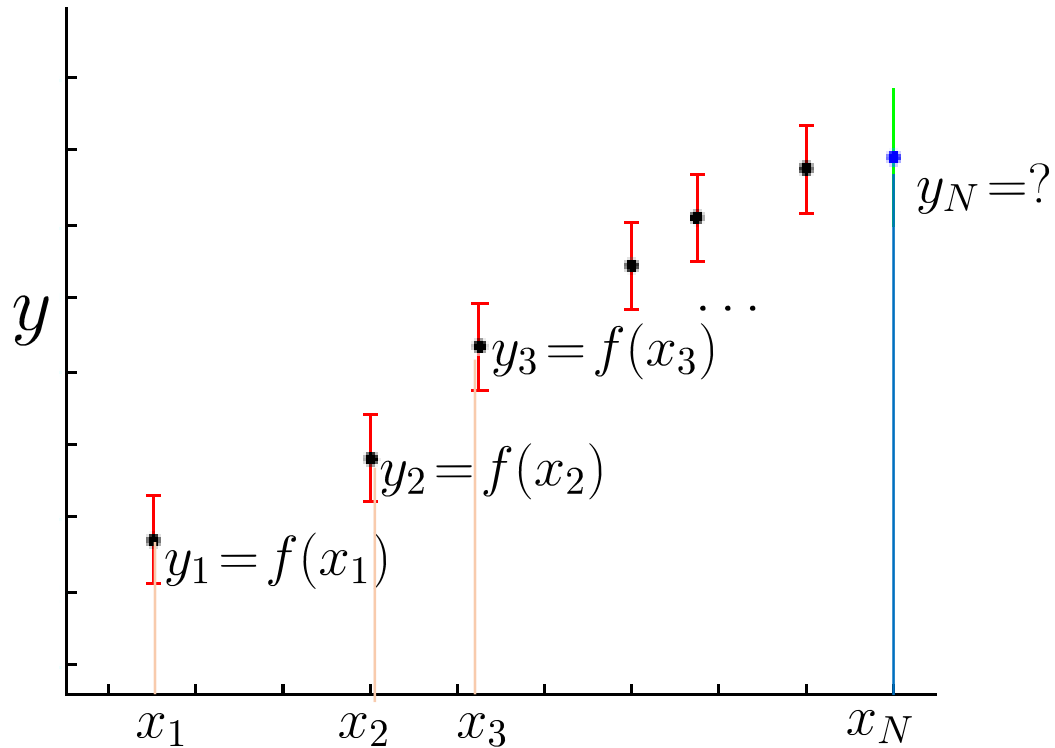


How to construct this joint distribution in GP?

By looking at the similarities in the continuous \mathbf{x} space, representing the locations at which we evaluate $y = f(\mathbf{x})$

Gaussian process - Interpretation

- A covariance over an arbitrarily large set of variables can be defined through the **covariance kernel function** $k(\mathbf{x}_i, \mathbf{x}_j)$, providing the covariance elements between any two sample locations \mathbf{x}_i and \mathbf{x}_j .



If x_N is similar to x_3 ,
we also expect y_N
to be similar to y_3 .

Gaussian process (GP)

For a set of input locations $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, a positive semidefinite covariance matrix (also known as the Gram matrix) is defined as

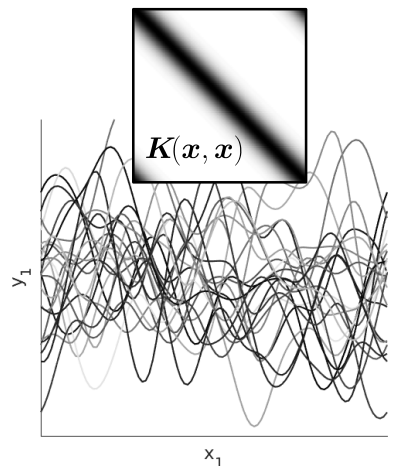
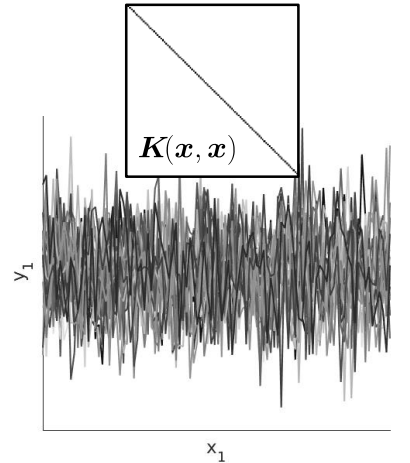
$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}.$$

The entire function evaluation $y_n = f(\mathbf{x}_n)$ associated with the set of inputs \mathbf{x}_n is a draw from a multivariate Gaussian distribution

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \mathbf{K}(\mathbf{x}, \mathbf{x})),$$

specifying a **distribution over functions**.

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \mathbf{K}(\mathbf{x}, \mathbf{x}))$$



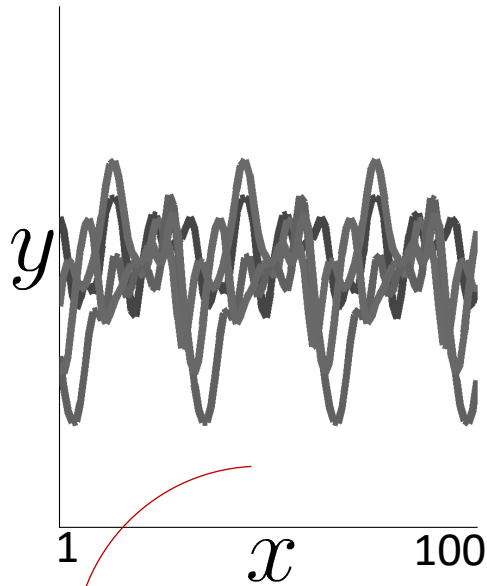
How to choose $k(\mathbf{x}_i, \mathbf{x}_j)$?

- We may know that our observations are samples from a process that is **smooth**, that is **continuous**, that has **typical amplitude**, or that the variations in the function take place within a **typical dynamic range**.
- These models require hyperparameters to be inferred, but **these hyperparameters define characteristics that are more generic** (such as the scale of a distribution) rather than acting explicitly on the structure or functional form of the signals.
- The notion of similarity will depend on the application: some of the basic aspects that can be defined through the covariance function $k(\mathbf{x}_i, \mathbf{x}_j)$ are the process **stationarity, isotropy, smoothness or periodicity**.
- With continuous time series, past observations will be informative about current data as a function of how long ago they were observed.
- This corresponds to a **stationary** covariance, dependent on the Euclidean distance $|\mathbf{x}_i - \mathbf{x}_j|$.
- This process is also considered as **isotropic** if it does not depend on directions between \mathbf{x}_i and \mathbf{x}_j .
- A process that is both stationary and isotropic is **homogeneous**.

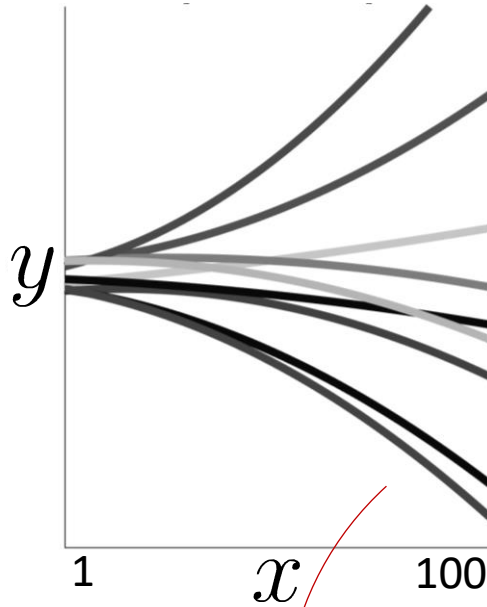
How to choose $k(\mathbf{x}_i, \mathbf{x}_j)$?

$$\mathbf{x} = [1, 2, \dots, 100]^\top$$

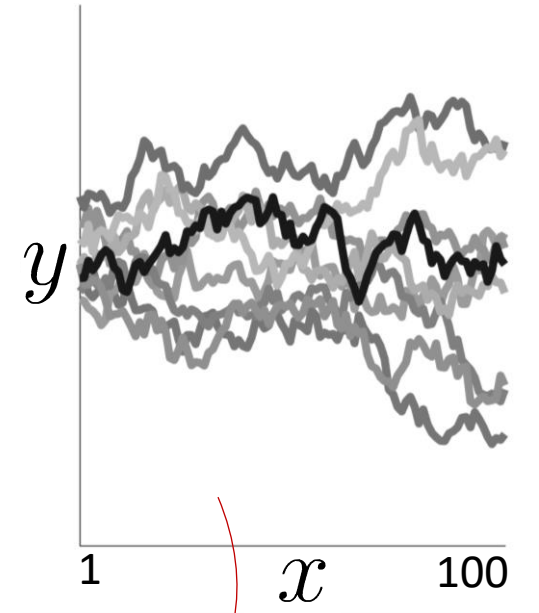
$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{x}, \mathbf{x}))$$



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}} \sin^2(\Theta_4^{\text{GP}} \|\mathbf{x}_i - \mathbf{x}_j\|)\right)$$



$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + \Theta_1^{\text{GP}})^2$$



$$k(\mathbf{x}_i, \mathbf{x}_j) = \min(\mathbf{x}_i, \mathbf{x}_j) + \Theta_1^{\text{GP}}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential function

A popular homogeneous covariance function is the **squared exponential kernel**, also known as **radial basis function**

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left(-\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right),$$

with two hyperparameters Θ_1^{GP} and Θ_2^{GP} corresponding respectively to **output and input scales** of the problem.

The radial basis function is widely employed when it is expected that nearby inputs \mathbf{x}_i and \mathbf{x}_j will have their corresponding outputs \mathbf{y}_i and \mathbf{y}_j also nearby (**assumption of continuity**).

Modeling noise in the observed y_n

If we assume there is noise associated with the observed function values $\mathbf{y}_n = f(\mathbf{x}_n) + \boldsymbol{\eta}$, this noise term can also be modeled in the covariance.

This noise is most often assumed to be uncorrelated from sample to sample, meaning that the noise term is only added to the diagonal elements of \mathbf{K} , giving a modified covariance for noisy observations of the form

$$\tilde{\mathbf{K}}(\mathbf{x}, \mathbf{x}) = \mathbf{K}(\mathbf{x}, \mathbf{x}) + \Theta_3^{\text{GP}} \mathbf{I}$$

where \mathbf{I} is the identity matrix and Θ_3^{GP} is a Gaussian process hyperparameter representing the noise variance.

Gaussian process regression (GPR)

We are interested in the **posterior distribution** of \mathbf{y}^* to be computed at some location(s) \mathbf{x}^* .

The **joint distribution** of the already observed \mathbf{y} (at location \mathbf{x}) augmented by \mathbf{y}^* (at location \mathbf{x}^*) is

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}^* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}(\mathbf{x}) \\ \boldsymbol{\mu}(\mathbf{x}^*) \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{x}, \mathbf{x}) & \mathbf{K}(\mathbf{x}, \mathbf{x}^*) \\ \mathbf{K}(\mathbf{x}^*, \mathbf{x}) & \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*) \end{bmatrix} \right)$$

$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

We can use the conditional probability property of Gaussians to evaluate the posterior distribution of \mathbf{y}^* , yielding a Gaussian

$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\text{with } \boldsymbol{\mu}^* = \boldsymbol{\mu}(\mathbf{x}^*) + \mathbf{K}(\mathbf{x}^*, \mathbf{x}) \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} (\mathbf{y} - \boldsymbol{\mu}(\mathbf{x}))$$

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

Gaussian process regression (GPR)

It is also often assumed in practice that $\begin{bmatrix} \boldsymbol{\mu}(\mathbf{x}) \\ \boldsymbol{\mu}(\mathbf{x}^*) \end{bmatrix} = \mathbf{0}$.

In this case, Gaussian processes can be completely defined by second-order statistics, where the Gram matrix \mathbf{K} is a positive semi-definite covariance.

Note that $\mathbf{K}(\mathbf{x}, \mathbf{x})^{-1}$ can be pre-computed so that the posterior distribution can be computed faster

$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\text{with } \boldsymbol{\mu}^* = \boldsymbol{\mu}(\mathbf{x}^*) + \mathbf{K}(\mathbf{x}^*, \mathbf{x}) \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} (\mathbf{y} - \boldsymbol{\mu}(\mathbf{x}))$$

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

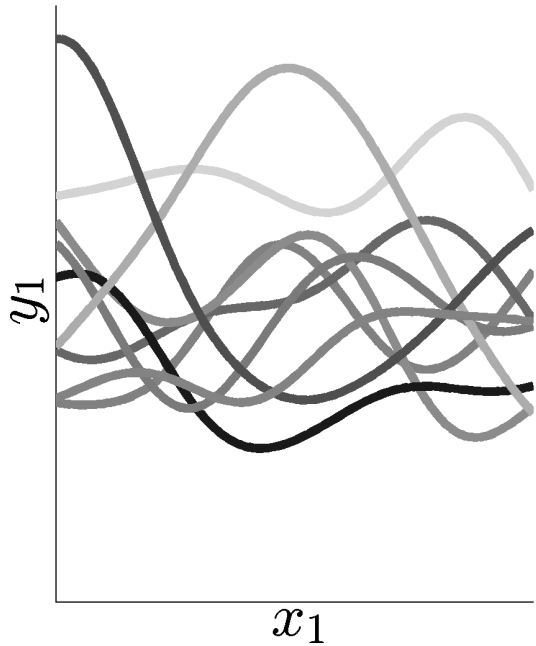
$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential function

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0$$

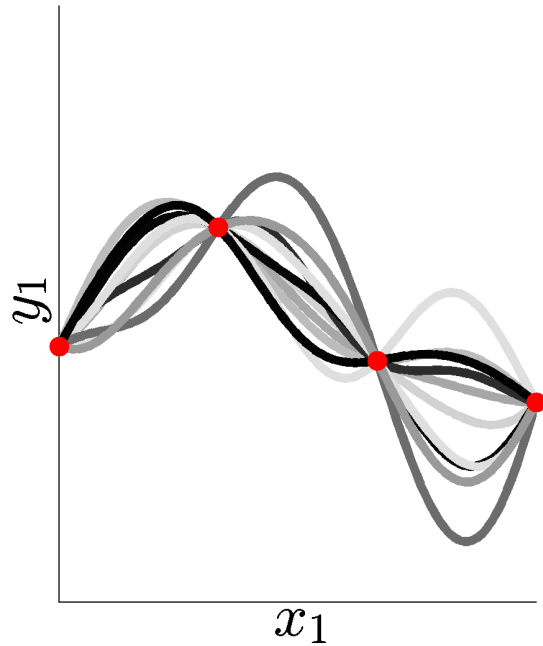
$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

Samples from prior



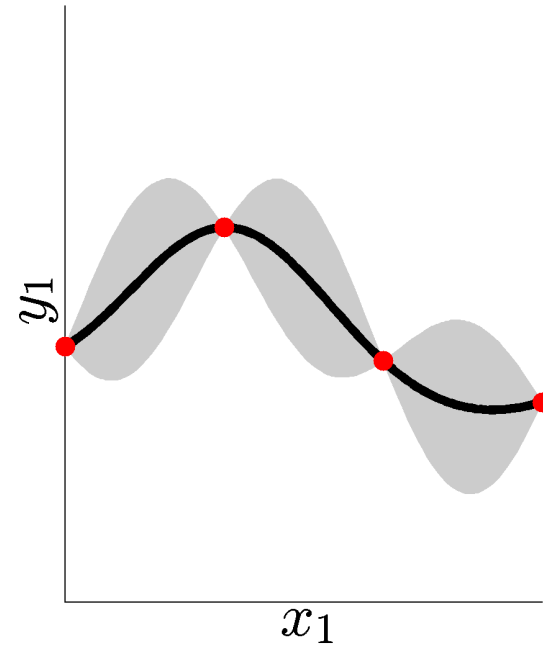
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Samples from posterior



$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Trajectory distribution



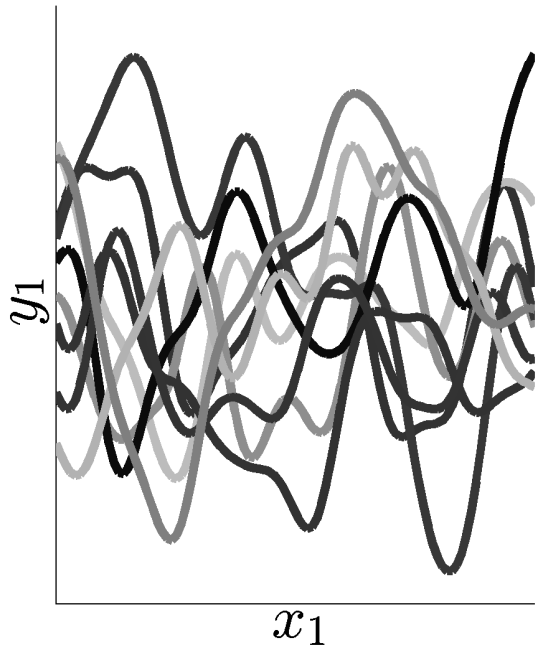
$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)\right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential function

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.01, \quad \Theta_3^{\text{GP}} = 0$$

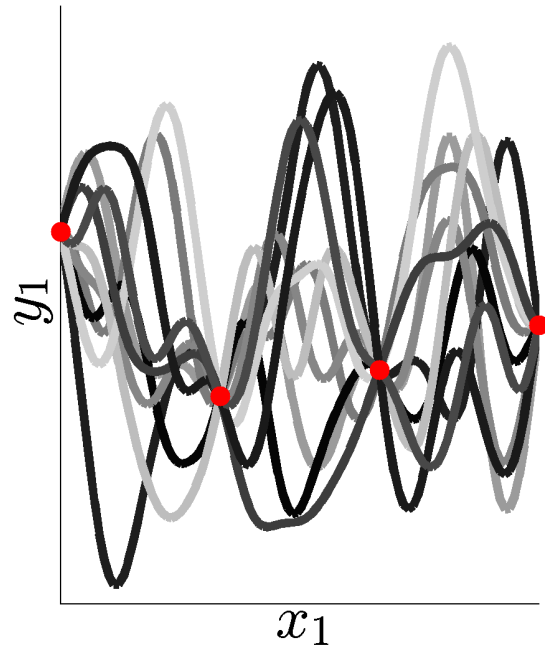
$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

Samples from prior



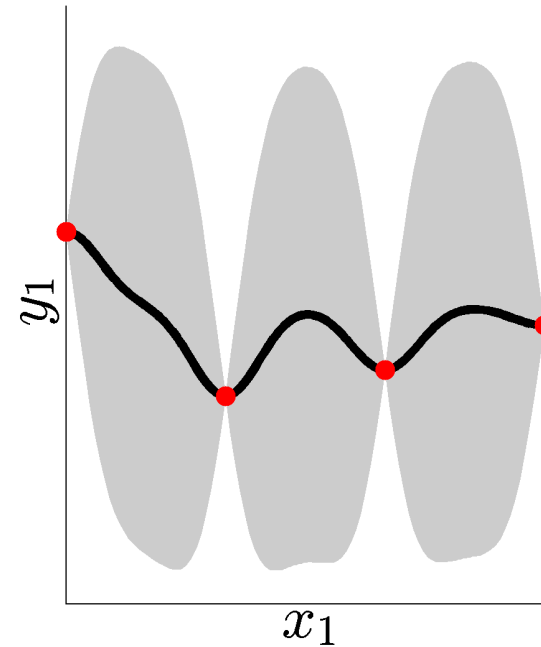
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Samples from posterior



$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Trajectory distribution



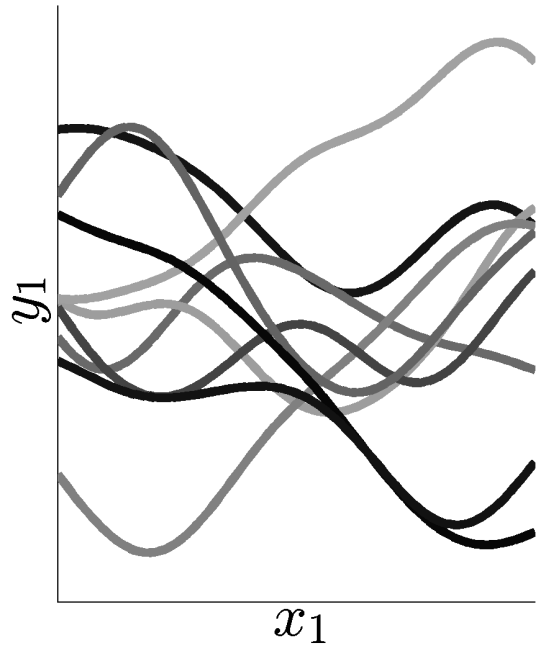
$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left(-\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential function

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0.01$$

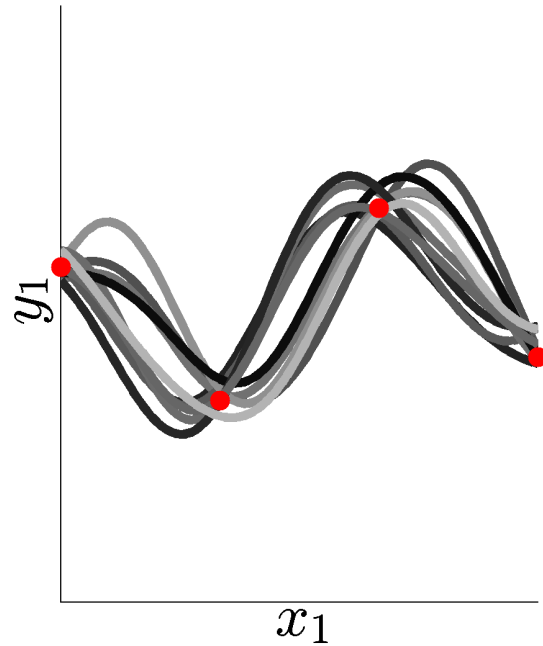
$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

Samples from prior



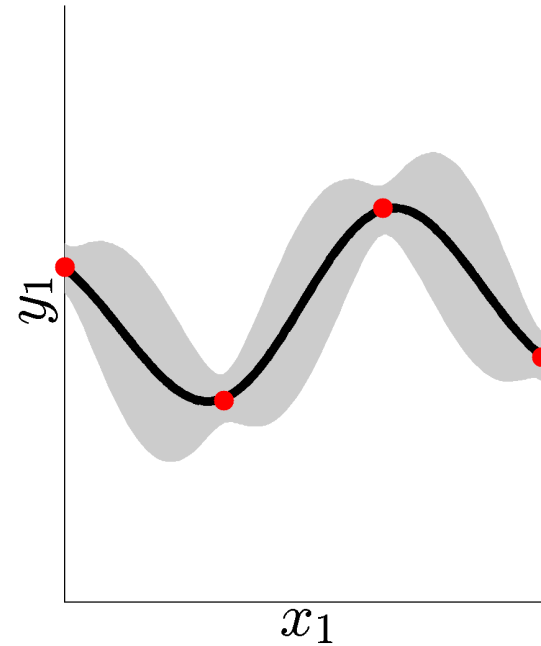
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Samples from posterior



$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}}(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)\right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential function

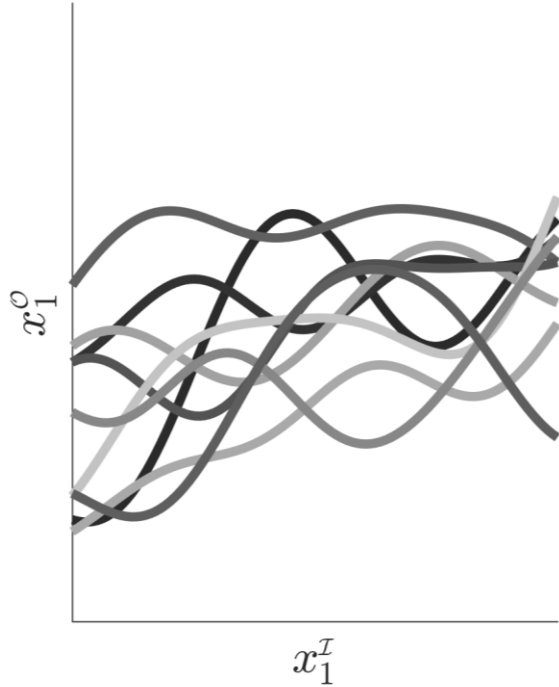
$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0.01$$

$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

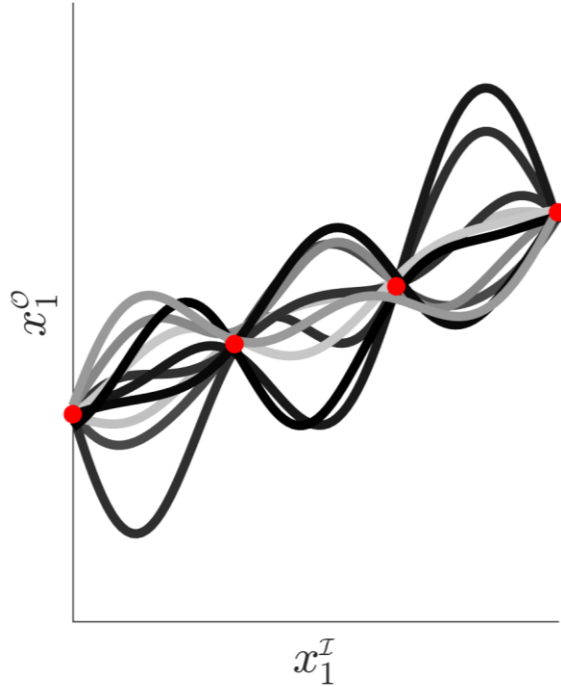
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

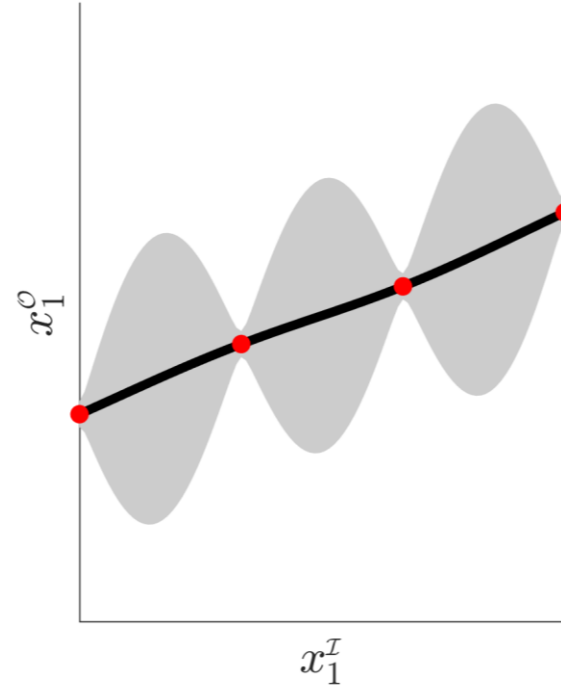
Samples from prior



Samples from posterior



Trajectory distribution



$$\boldsymbol{\mu}(\mathbf{x}) = \alpha \mathbf{x}$$

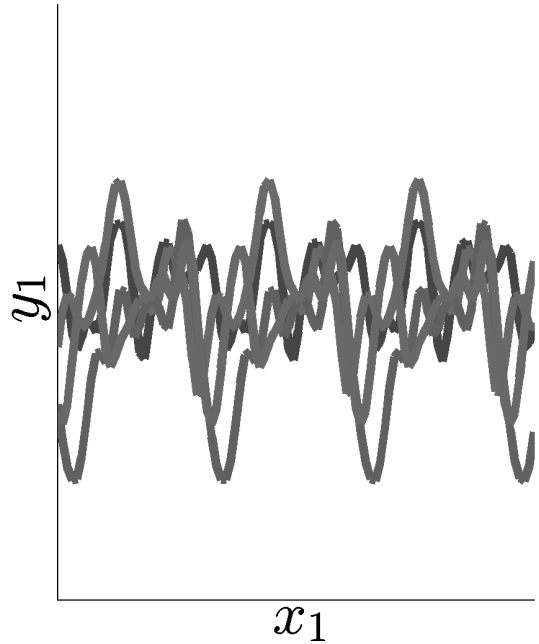
$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)\right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as **periodic** covariance function

$$\Theta_1^{\text{GP}} = 0.1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0, \quad \Theta_4^{\text{GP}} = 10$$

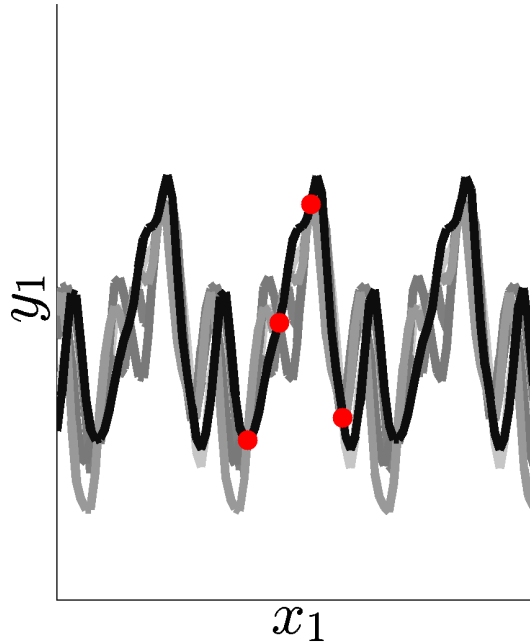
$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

Samples from prior



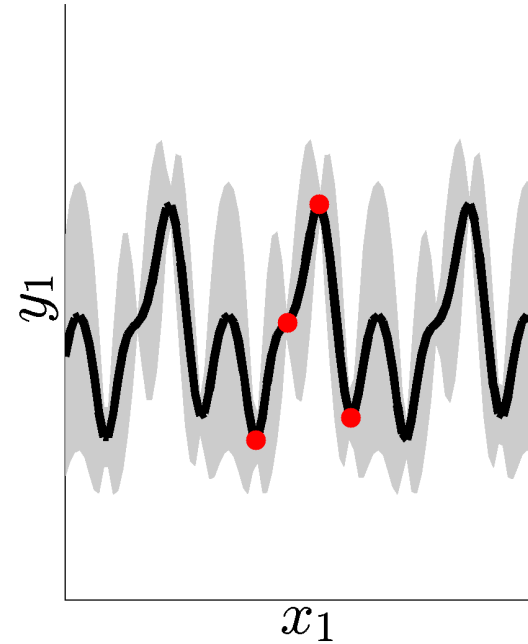
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Samples from posterior

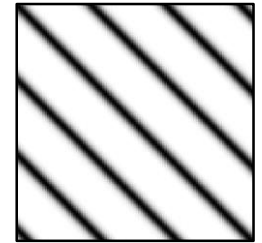


$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

Trajectory distribution



$$\mathbf{K}(\mathbf{x}^*, \mathbf{x}^*)$$



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left(-\frac{1}{\Theta_2^{\text{GP}}} \sin^2 \left(\Theta_4^{\text{GP}} \|\mathbf{x}_i - \mathbf{x}_j\| \right) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as Matérn covariance function

Another popular covariance kernel function is the Matérn function

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} \frac{d}{\rho} \right)^\nu K_\nu \left(\sqrt{2\nu} \frac{d}{\rho} \right)$$

$$\text{with } d = \|\mathbf{x}_i - \mathbf{x}_j\|$$

where Γ is the gamma function, K_ν is the modified Bessel function of the second kind, and ρ and ν are non-negative parameters of the covariance.

A Gaussian process with Matérn covariance has sample paths that are $\lfloor \nu - 1 \rfloor$ times differentiable.

$k(x_i, x_j)$ as Matérn covariance function

Simplification for ν half integer

When $\nu = p + 1/2$, $p \in \mathbb{N}^+$, the Matérn covariance can be written as a product of an exponential and a polynomial of order p :

$$C_{p+1/2}(d) = \sigma^2 \exp\left(-\frac{\sqrt{2\nu}d}{\rho}\right) \frac{\Gamma(p+1)}{\Gamma(2p+1)} \sum_{i=0}^p \frac{(p+i)!}{i!(p-i)!} \left(\frac{\sqrt{8\nu}d}{\rho}\right)^{p-i}$$

For $\nu = 1/2$ ($p = 0$): $C_{1/2}(d) = \sigma^2 \exp\left(-\frac{d}{\rho}\right)$

For $\nu = 3/2$ ($p = 1$): $C_{3/2}(d) = \sigma^2 \left(1 + \frac{\sqrt{3}d}{\rho}\right) \exp\left(-\frac{\sqrt{3}d}{\rho}\right)$

For $\nu = 5/2$ ($p = 2$): $C_{5/2}(d) = \sigma^2 \left(1 + \frac{\sqrt{5}d}{\rho} + \frac{5d^2}{3\rho^2}\right) \exp\left(-\frac{\sqrt{5}d}{\rho}\right)$

As $\nu \rightarrow \infty$, the Matérn covariance converges to the squared exponential covariance function.

$k(\mathbf{x}_i, \mathbf{x}_j)$ as **Matérn** covariance function

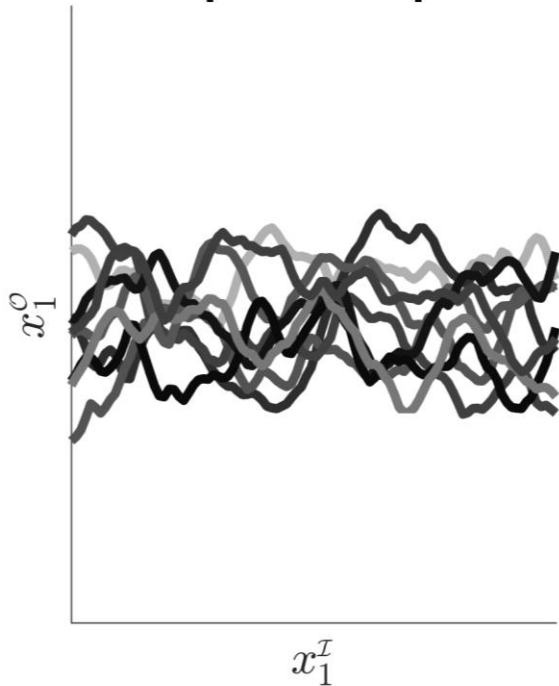
$$\Theta_1^{\text{GP}} = 0.1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0.0001$$

$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

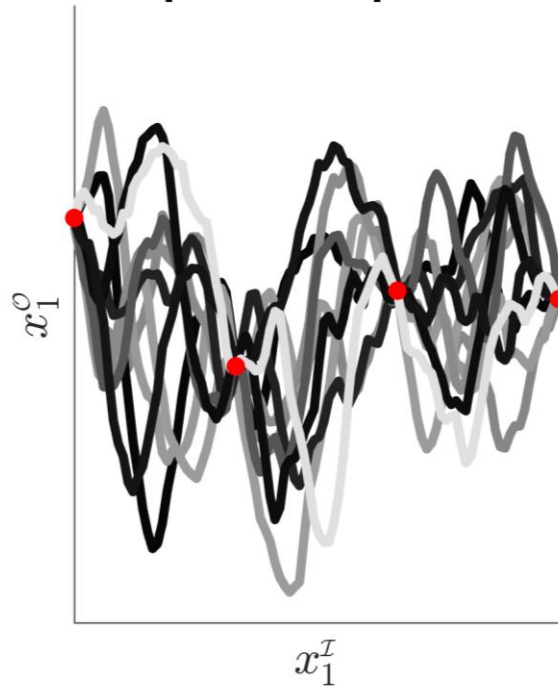
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

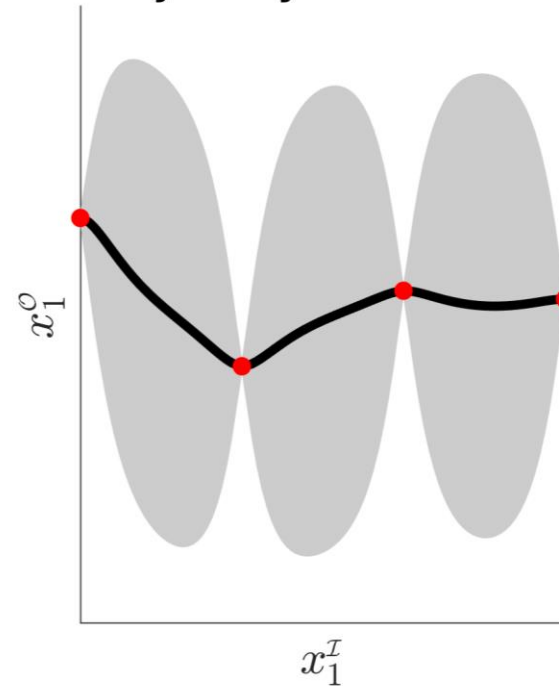
Samples from prior



Samples from posterior



Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \left(1 + \frac{\sqrt{3} \|\mathbf{x}_i - \mathbf{x}_j\|}{\Theta_2^{\text{GP}}} \right) \exp \left(-\frac{\sqrt{3} \|\mathbf{x}_i - \mathbf{x}_j\|}{\Theta_2^{\text{GP}}} \right)$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as **Brownian motion** covariance function

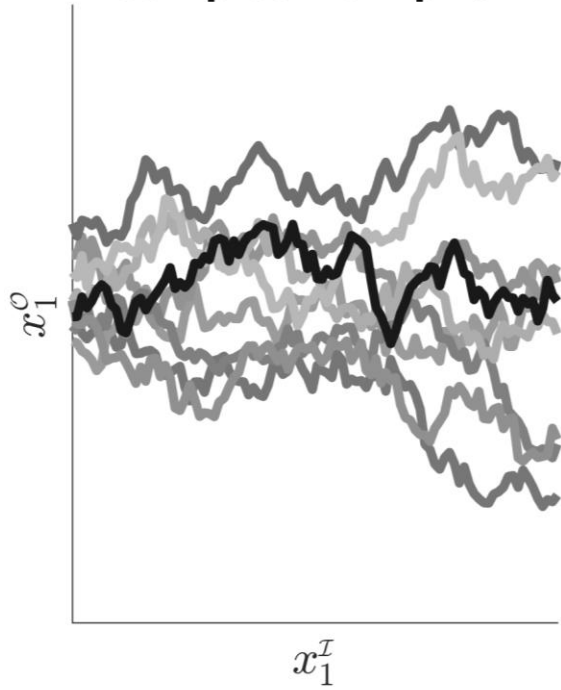
The **Wiener process** is a simple continuous-time stochastic process often put in connection to the Brownian motion.

$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

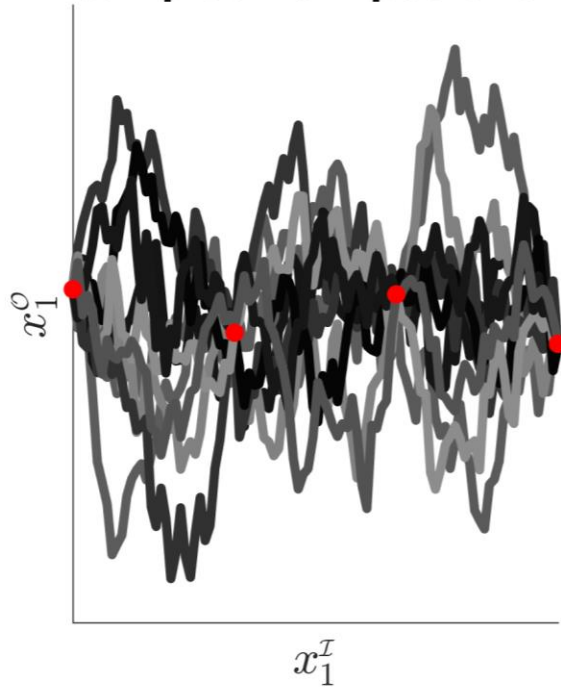
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

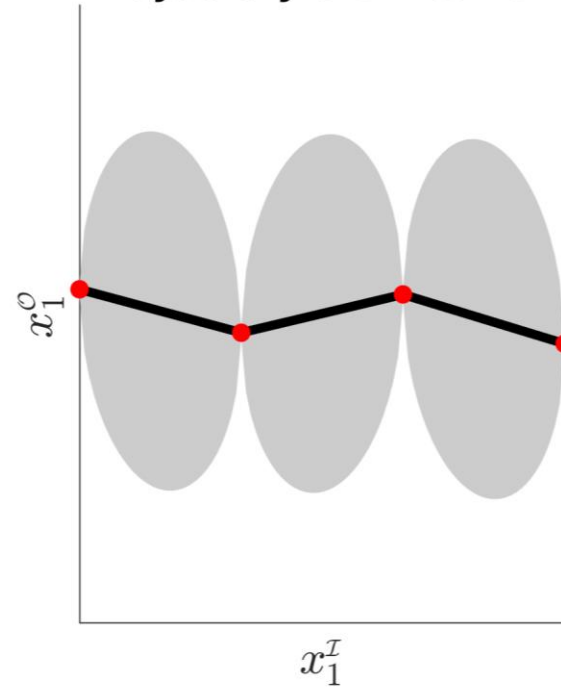
Samples from prior



Samples from posterior



Trajectory distribution



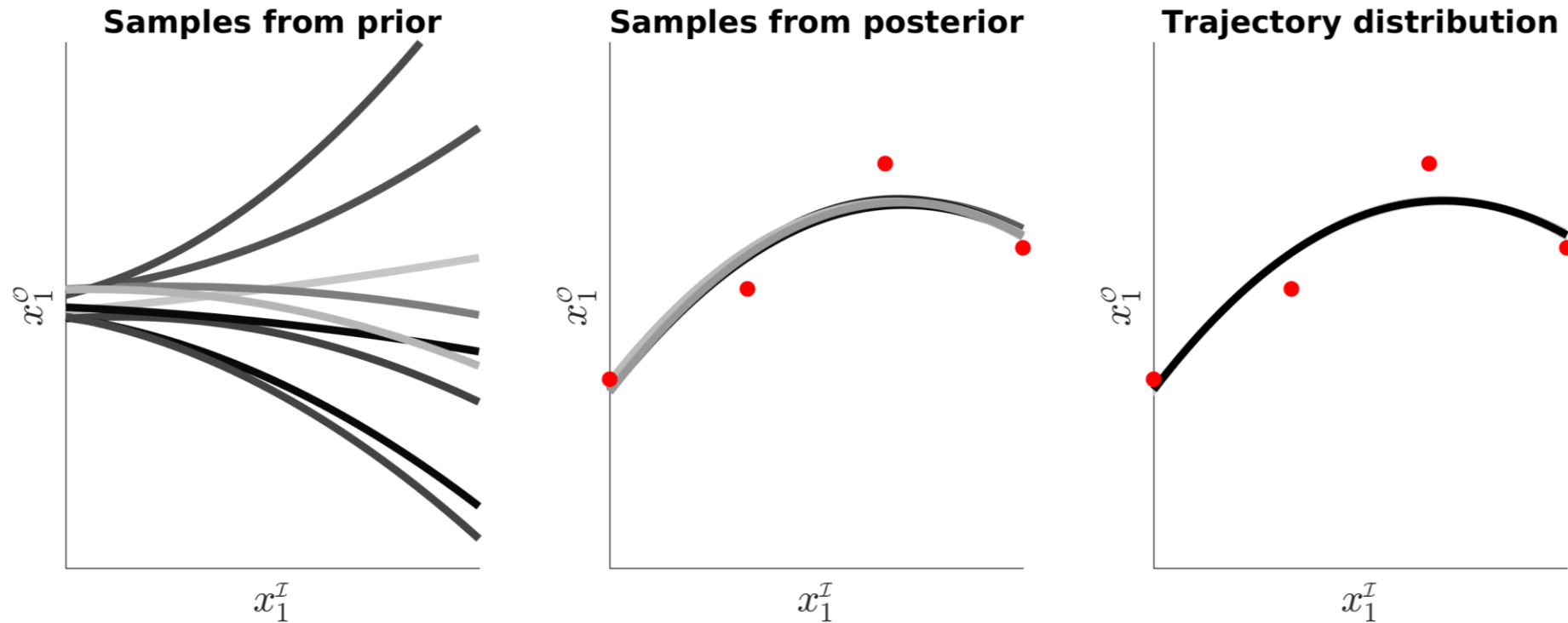
$$k(\mathbf{x}_i, \mathbf{x}_j) = \min(\mathbf{x}_i, \mathbf{x}_j) + \Theta_1^{\text{GP}}$$

$$\Theta_1^{\text{GP}} = 0.1$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as **quadratic** covariance function

Bayesian linear regression is equivalent to a GP with covariance function $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$.

$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*)) \quad \mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*) \quad \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$



$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + \Theta_1^{\text{GP}})^2 \quad \Theta_1^{\text{GP}} = 0.1$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as **polynomial** covariance function

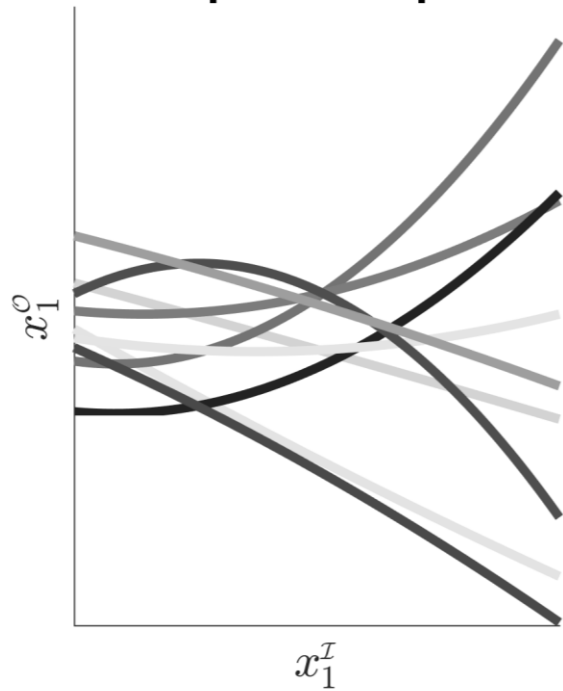
$$\Theta_1^{\text{GP}} = 0.1$$

$$\mathbf{y}^* \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

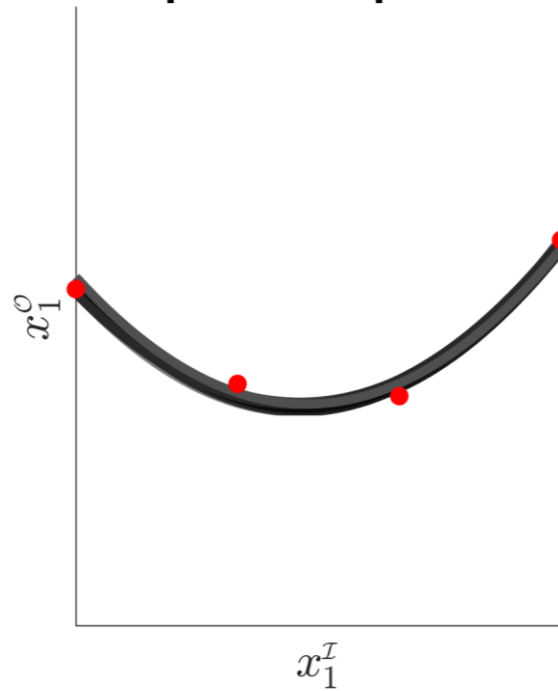
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

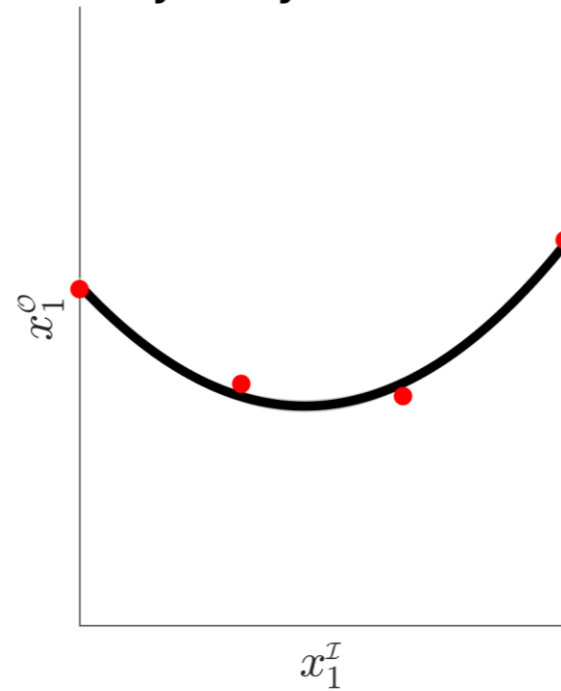
Samples from prior



Samples from posterior

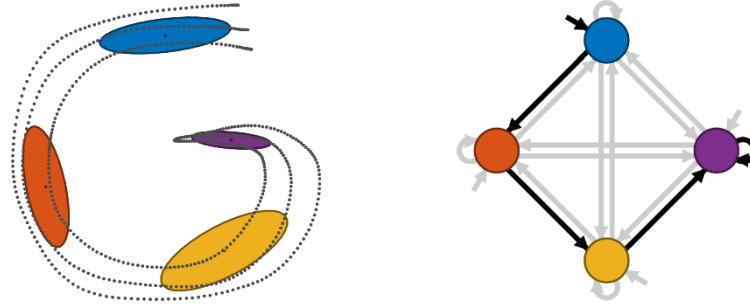


Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j)^2 + \mathbf{x}_i^\top \mathbf{x}_j + \Theta_1^{\text{GP}}$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ as probabilistic model covariance



- Another powerful approach to the construction of kernels is to exploit probabilistic models.
- Given a generative model $P(\mathbf{x})$, a kernel can be defined as $k(\mathbf{x}_i, \mathbf{x}_j) = P(\mathbf{x}_i) P(\mathbf{x}_j)$, which can be interpreted as an inner product in the one-dimensional feature space defined by the mapping $P(\mathbf{x})$.
- Namely, two inputs \mathbf{x}_i and \mathbf{x}_j will be similar if they both have high probabilities to belong to the model.
- This can bring additional properties to the underlying process such as the capability of handling missing data or partial sequences of various lengths (e.g., with HMM).

$k(\mathbf{x}_i, \mathbf{x}_j)$ as weighted sum of covariance functions

- A covariance function can be defined as a **linear combination of other covariance functions**, which can be exploited to incorporate **different insights about the dataset**.
- Such an approach can be exploited as an alternative to optimizing kernel parameters (also known as multiple kernel learning).
- The idea is to define the kernel as a **weighted sum of basis kernels**, and then to **optimize the weights instead of the kernel parameters**.

$$\text{Dictionary of basis kernel functions} \left\{ \begin{array}{l} k_1(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j)^2 + \mathbf{x}_i^\top \mathbf{x}_j \\ k_2(\mathbf{x}_i, \mathbf{x}_j) = \min(\mathbf{x}_i, \mathbf{x}_j) \\ k_3(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)\right) \end{array} \right.$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} k_1(\mathbf{x}_i, \mathbf{x}_j) + \Theta_2^{\text{GP}} k_2(\mathbf{x}_i, \mathbf{x}_j) + \Theta_3^{\text{GP}} k_3(\mathbf{x}_i, \mathbf{x}_j)$$

Some extensions of Gaussian processes

- **Cokriging:**

Extending GPR to multiple target variables y .

- **Sparse GP:**

A known bottleneck in Gaussian process prediction is that the computational complexity of prediction is $O(N^3)$

→ **not feasible for large data sets!**

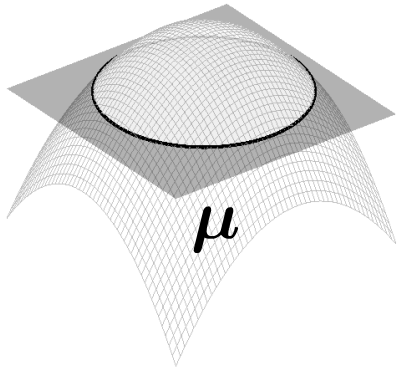
Sparse Gaussian processes circumvent this issue by building a representative set for the given process $y = f(\mathbf{x})$.

- **Gaussian process latent variable models (GPLVM):**

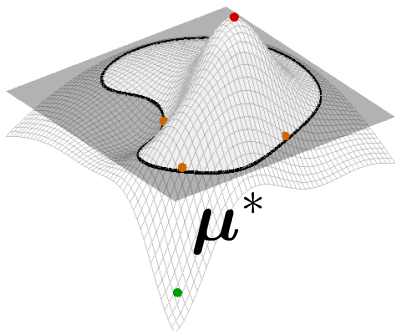
GPLVM is a probabilistic dimensionality reduction method that uses GPs to find a lower dimensional non-linear embedding of high dimensional data.

Example: Gaussian Process Implicit Surface (GPIS)

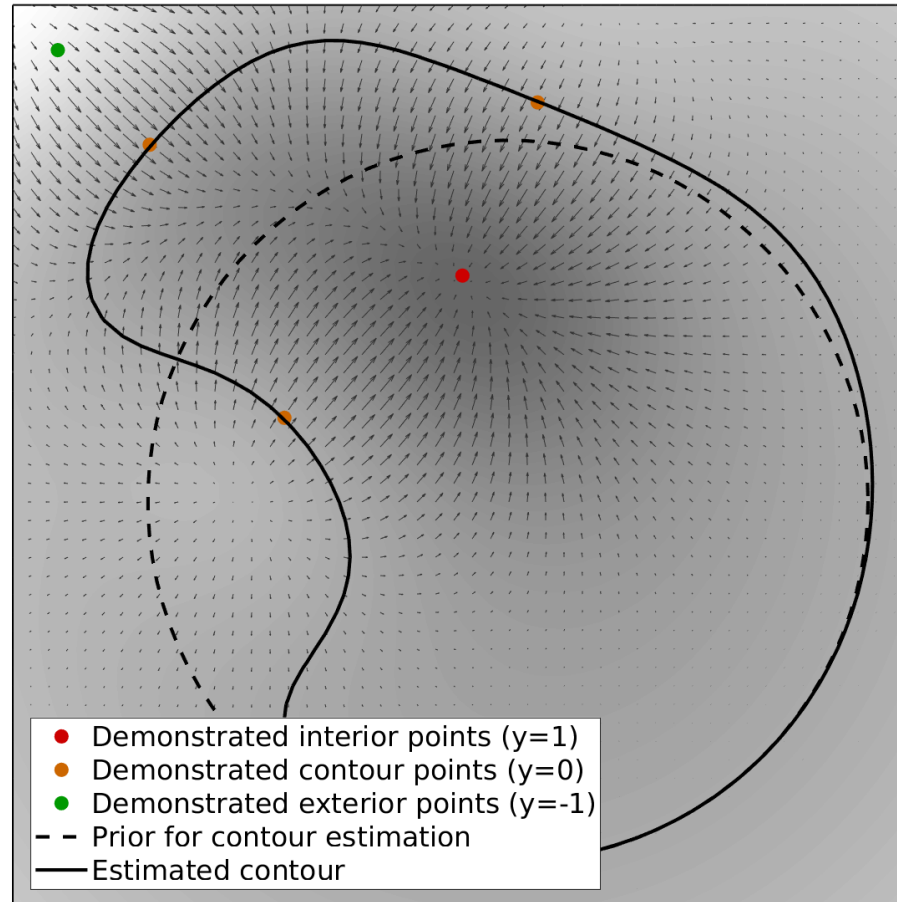
Prior (circular contour)



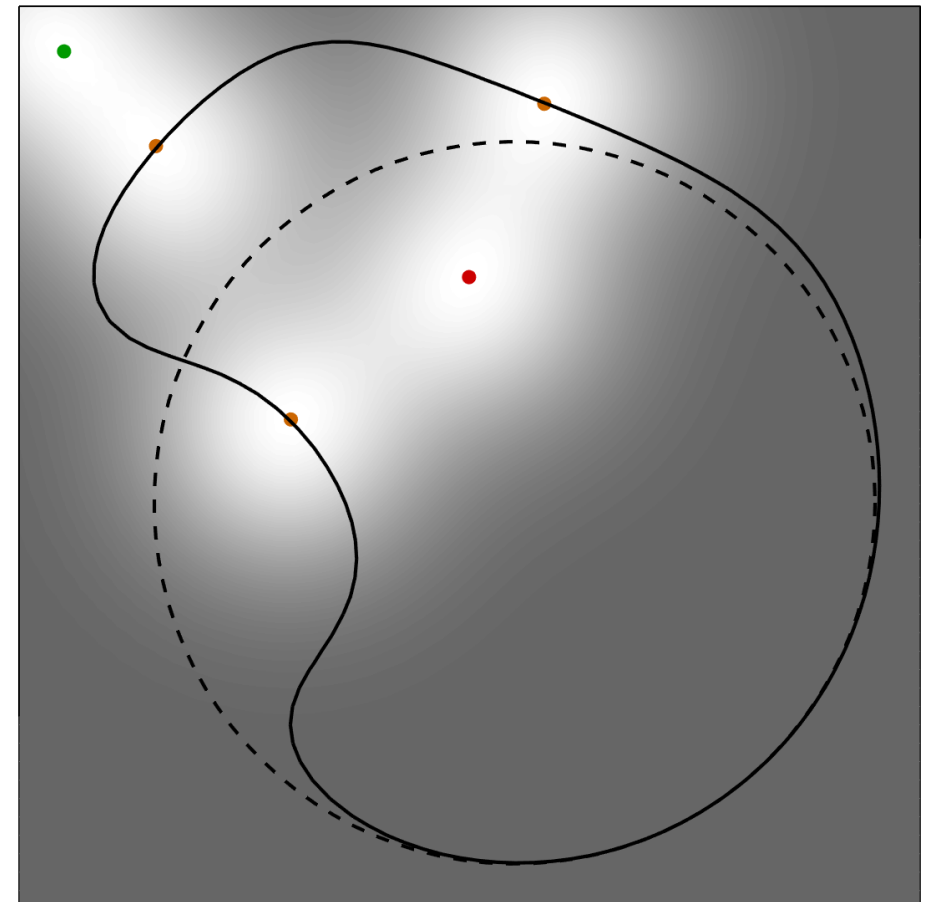
Posterior



Distance to contour and gradient



Uncertainty



Summary of today's lecture

Properties of multivariate Gaussian distributions:

- Product of Gaussians
- Linear transformation and combination
- Conditional distribution

Three nonlinear regression models:

- Locally weighted regression (LWR)
- Gaussian mixture regression (GMR)
- Gaussian process regression (GPR)

**Modeling possible
co-variations**
(a.k.a. aleatoric uncertainty)

**Modeling uncertainty
of the estimate**
(a.k.a. epistemic uncertainty)

References

LWR

C. G. Atkeson, A. W. Moore, and S. Schaal. Locally weighted learning for control. *Artificial Intelligence Review*, 11(1-5):75–113, 1997

W.S. Cleveland. Robust locally weighted regression and smoothing scatterplots. *American Statistical Association* 74(368):829–836, 1979

GMR

Z. Ghahramani and M. I. Jordan. Supervised learning from incomplete data via an EM approach. In *Advances in Neural Information Processing Systems (NIPS)*, volume 6, pages 120–127, 1994

S. Calinon. *Mixture models for the analysis, edition, and synthesis of continuous time series*. Mixture Models and Applications, Springer, 2019

GPR

C.K.I. Williams and C.E. Rasmussen. Gaussian processes for regression. In *Advances in Neural Information Processing Systems (NIPS)*, pages 514–520, 1996

C.E. Rasmussen and C.K.I. Williams. *Gaussian processes for machine learning*. MIT Press, Cambridge, MA, USA, 2006

S. Roberts, M. Osborne, M. Ebdon, S. Reece, N. Gibson, and S. Aigrain. Gaussian processes for time-series modelling. *Philosophical Trans. of the Royal Society A*, 371(1984):1–25, 2012

GPIS

O. Williams and A. Fitzgibbon. Gaussian Process Implicit Surfaces. In *Gaussian Processes in Practice*, 2007