

EE613
Machine Learning for Engineers

NONLINEAR REGRESSION I

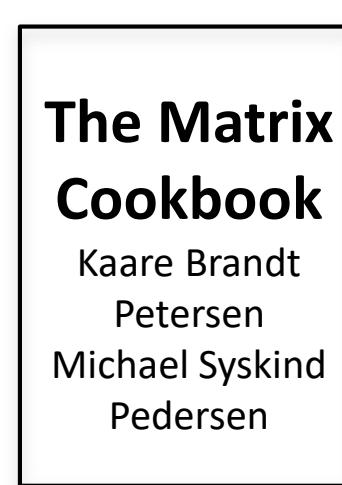
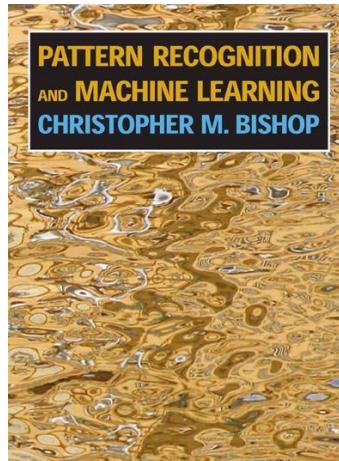
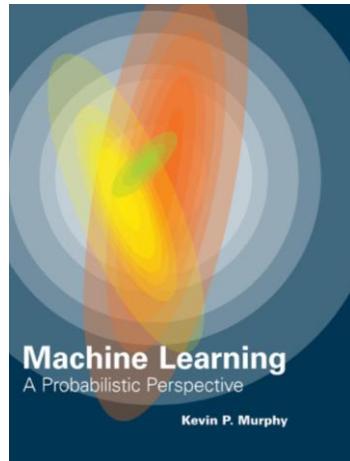
Sylvain Calinon
Robot Learning & Interaction Group
Idiap Research Institute
Dec. 13, 2017

Outline

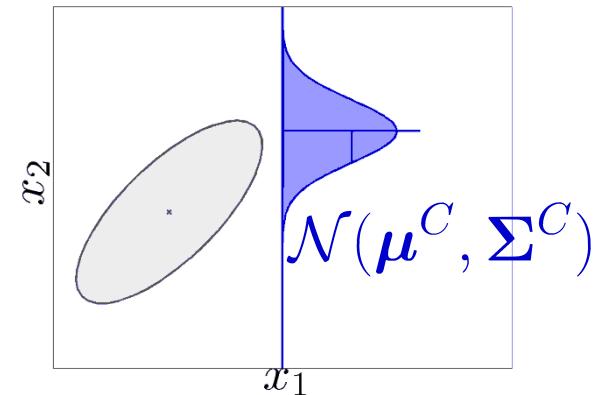
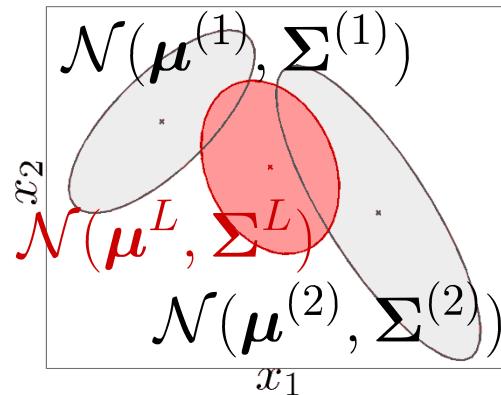
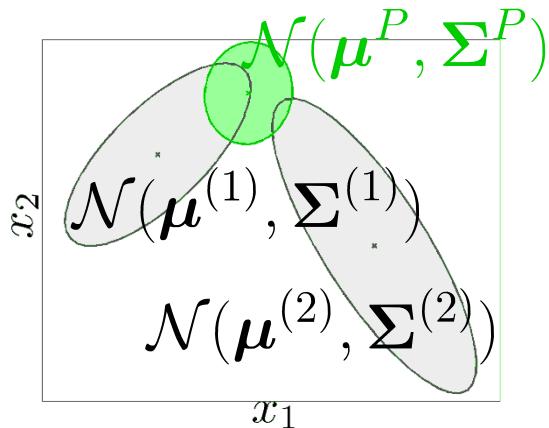
- Properties of multivariate Gaussian distributions
Matlab codes: *demo_Gaussian_product01.m*,
demo_Gaussian_conditioning01.m,
demo_Gaussian_conditioning_noisyInput01.m,
demo_Gaussian_lawTotalCov01.m
- Locally weighted regression (LWR)
Matlab code: *demo_LWR01.m*,
- Dynamical movement primitives (DMP)
Matlab code: *demo_DMP01.m*
- Gaussian mixture regression (GMR)
Matlab codes: *demo_GMR01.m*, *demo_GMR_polyFit01.m*,
demo_DMP_GMR01.m

Properties of multivariate Gaussian distributions

Matlab codes: `demo_Gaussian_product01.m`,
`demo_Gaussian_conditioning01.m`,
`demo_Gaussian_conditioning_noisyInput01.m`,
`demo_Gaussian_lawTotalCov01.m`



Some very useful properties



Product of Gaussians:

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

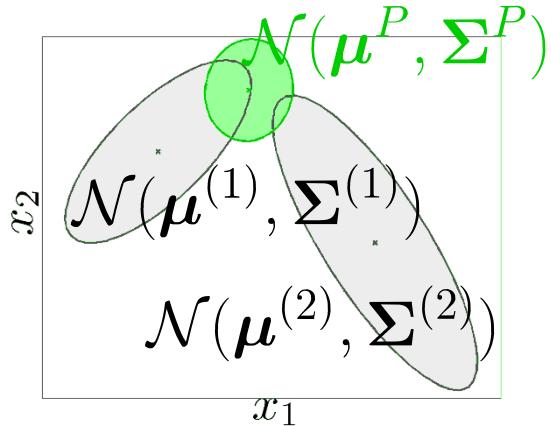
Linear combination:

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

Conditional probability:

$$\mathcal{N}(\mu^C, \Sigma^C) \sim \mathcal{P}(x_2|x_1)$$

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

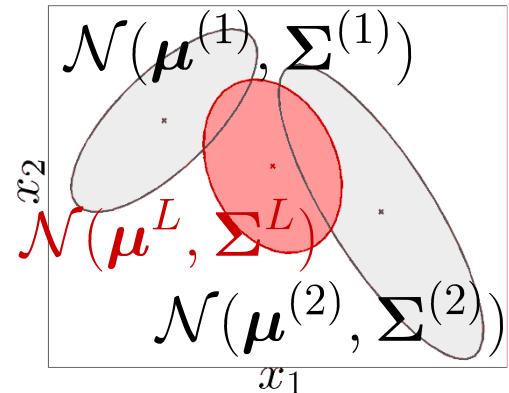
$$c \mathcal{N}(\boldsymbol{\mu}^P, \boldsymbol{\Sigma}^P) = \mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)}) \cdot \mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)}),$$

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

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

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

Linear combination



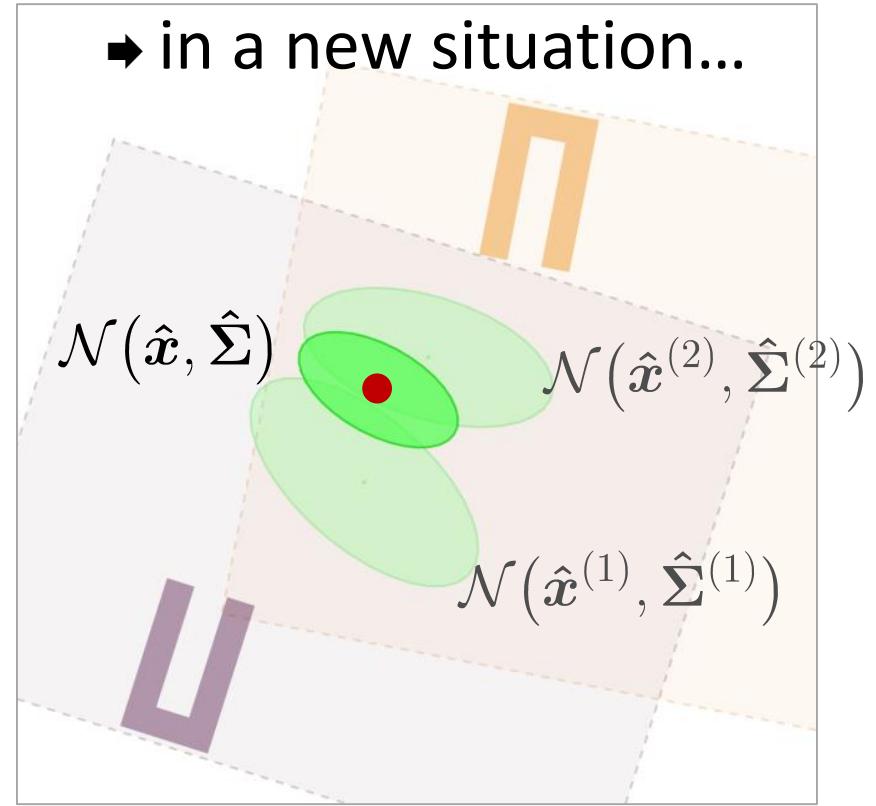
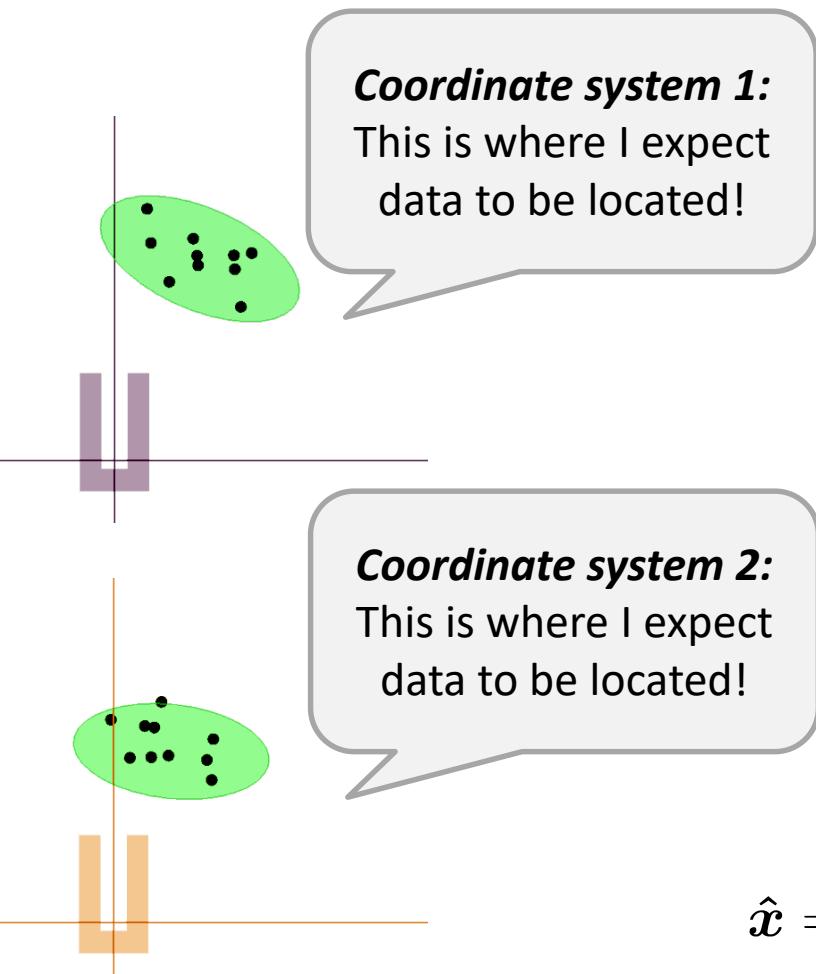
If $\mathbf{x}^{(1)} \sim \mathcal{N}(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)})$ and $\mathbf{x}^{(2)} \sim \mathcal{N}(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)})$, the linear transformation $\mathbf{A}^{(1)}\mathbf{x}^{(1)} + \mathbf{A}^{(2)}\mathbf{x}^{(2)} + \mathbf{c}$ follows the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

$$\begin{aligned}\boldsymbol{\mu}^L &= \mathbf{A}^{(1)}\boldsymbol{\mu}^{(1)} + \mathbf{A}^{(2)}\boldsymbol{\mu}^{(2)} + \mathbf{c}, \\ \boldsymbol{\Sigma}^L &= \mathbf{A}^{(1)}\boldsymbol{\Sigma}^{(1)}\mathbf{A}^{(1)\top} + \mathbf{A}^{(2)}\boldsymbol{\Sigma}^{(2)}\mathbf{A}^{(2)\top}.\end{aligned}$$

Example: Fusion of sensor/control information

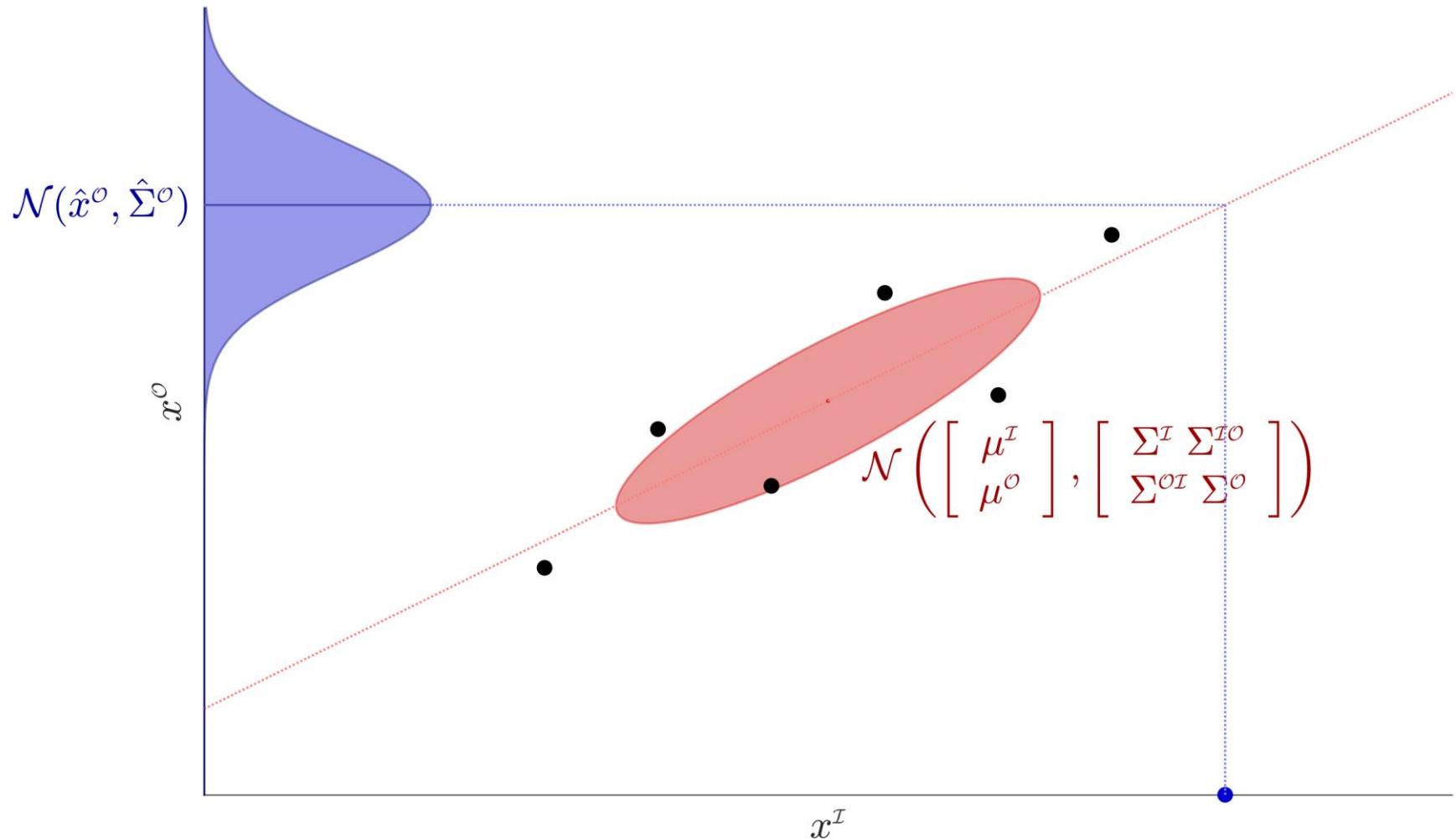


$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \sum_{j=1}^2 (\mathbf{x} - \hat{\mathbf{x}}^{(j)})^\top \hat{\Sigma}^{(j)-1} (\mathbf{x} - \hat{\mathbf{x}}^{(j)})$$

→ Product of linearly transformed Gaussians

Conditional probability

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



→ Linear regression from joint distribution

Conditional probability

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

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

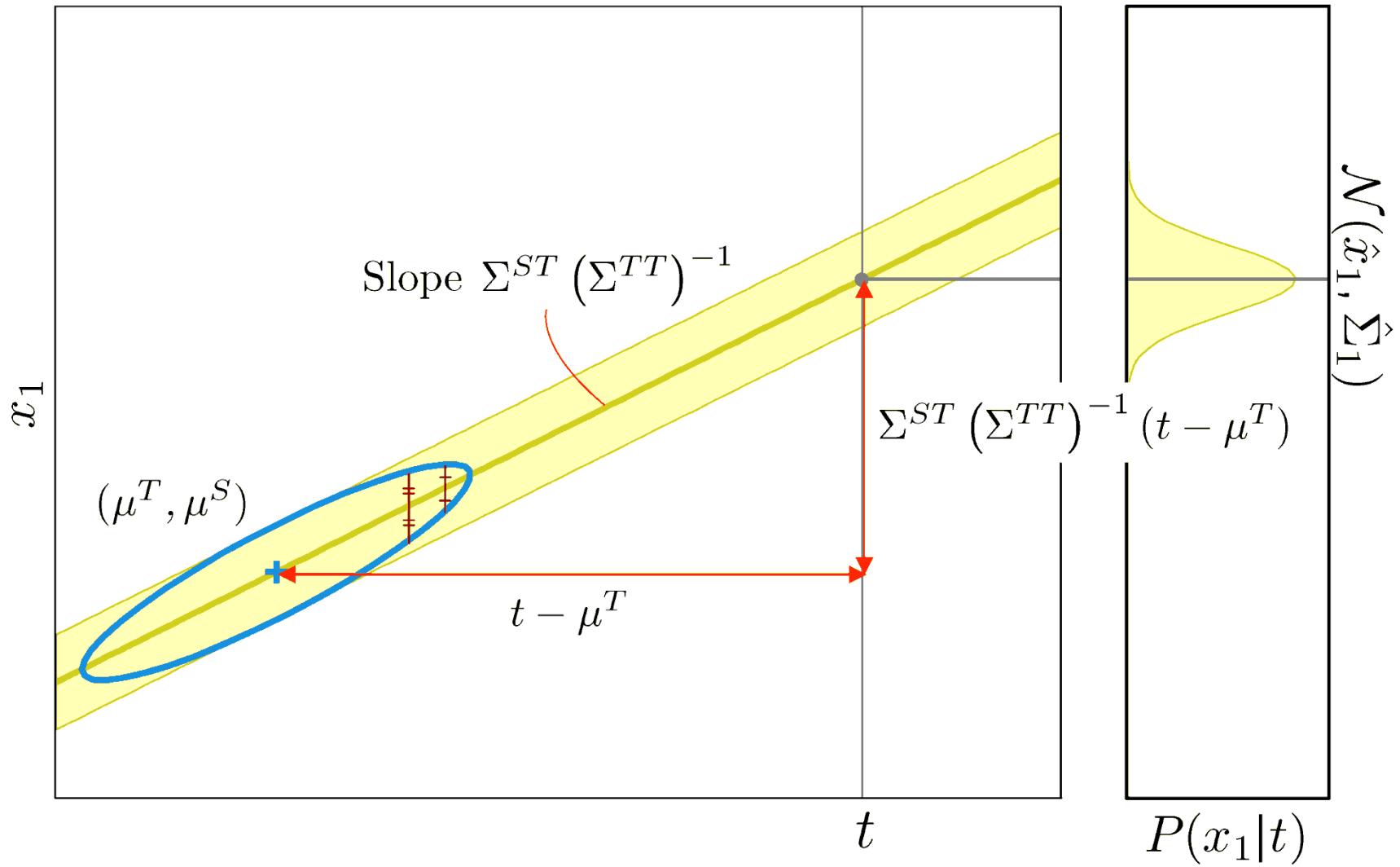
If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we have that $\mathbf{x}^{\mathcal{O}} | \mathbf{x}^{\mathcal{I}} \sim \mathcal{N}(\hat{\mathbf{x}}^{\mathcal{O}}, \hat{\boldsymbol{\Sigma}}^{\mathcal{O}})$, with parameters

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

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

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

Conditional probability - Geometric interpretation



Conditional probability - Proof

$$\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}^{\mathcal{I}} \\ \boldsymbol{x}^{\mathcal{O}} \end{bmatrix} \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}^{\mathcal{I}} \\ \boldsymbol{\mu}^{\mathcal{O}} \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}^{\mathcal{I}} & \boldsymbol{\Sigma}^{\mathcal{I}\mathcal{O}} \\ \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} & \boldsymbol{\Sigma}^{\mathcal{O}} \end{bmatrix}$$

We want to find the distribution of $\boldsymbol{x}^{\mathcal{O}}$ that maximizes the log-likelihood

$$\begin{aligned} f(\boldsymbol{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \log (\mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})) \\ &= -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) - \frac{D}{2} \log(2\pi), \end{aligned}$$

when $\boldsymbol{x}^{\mathcal{I}}$ is known and acts as a constant.

This can be computed by deriving the above equation and equating to zero, namely

$$\frac{\partial f}{\partial \boldsymbol{x}^{\mathcal{O}}} = 0.$$

Conditional probability - Proof

$$\Sigma = \begin{bmatrix} \Sigma^{\mathcal{I}} & \Sigma^{\mathcal{I}\mathcal{O}} \\ \Sigma^{\mathcal{O}\mathcal{I}} & \Sigma^{\mathcal{O}} \end{bmatrix}$$

To do this, we first note that Σ^{-1} can be partitioned as

$$\begin{aligned} \Sigma^{-1} = \Gamma &= \begin{bmatrix} \Gamma^{\mathcal{I}} & \Gamma^{\mathcal{I}\mathcal{O}} \\ \Gamma^{\mathcal{O}\mathcal{I}} & \Gamma^{\mathcal{O}} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & -\Sigma^{\mathcal{I}-1}\Sigma^{\mathcal{I}\mathcal{O}} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \Sigma^{\mathcal{I}-1} & \mathbf{0} \\ \mathbf{0} & S^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\Sigma^{\mathcal{O}\mathcal{I}}\Sigma^{\mathcal{I}-1} & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \Sigma^{\mathcal{I}-1} + \Sigma^{\mathcal{I}-1}\Sigma^{\mathcal{I}\mathcal{O}}S^{-1}\Sigma^{\mathcal{O}\mathcal{I}}\Sigma^{\mathcal{I}-1} & -\Sigma^{\mathcal{I}-1}\Sigma^{\mathcal{I}\mathcal{O}}S^{-1} \\ -S^{-1}\Sigma^{\mathcal{O}\mathcal{I}}\Sigma^{\mathcal{I}-1} & S^{-1} \end{bmatrix}, \end{aligned}$$

where $S = \Sigma^{\mathcal{O}} - \Sigma^{\mathcal{O}\mathcal{I}}\Sigma^{\mathcal{I}-1}\Sigma^{\mathcal{I}\mathcal{O}}$ is the **Schur complement** of Σ .

The above result can be shown by using a LDU decomposition of Σ , where D is a diagonal matrix and L and U are atomic triangular matrices (lower and upper, respectively), and then computing its inverse by exploiting the inversion properties of diagonal and atomic triangular matrices.

Conditional probability - Proof

$$\Gamma = \begin{bmatrix} \Gamma^{\mathcal{I}} & \Gamma^{\mathcal{I}\mathcal{O}} \\ \Gamma^{\mathcal{O}\mathcal{I}} & \Gamma^{\mathcal{O}} \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma^{\mathcal{I}} & \Sigma^{\mathcal{I}\mathcal{O}} \\ \Sigma^{\mathcal{O}\mathcal{I}} & \Sigma^{\mathcal{O}} \end{bmatrix}$$

With this partitioning, we can see that

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Gamma (\mathbf{x} - \boldsymbol{\mu}) &= -\frac{1}{2}(\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}})^\top \Gamma^{\mathcal{I}} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}}) \\ &\quad - \frac{1}{2}(\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}})^\top \Gamma^{\mathcal{I}\mathcal{O}} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}}) \\ &\quad - \frac{1}{2}(\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}})^\top \Gamma^{\mathcal{O}\mathcal{I}} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}}) \\ &\quad - \frac{1}{2}(\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}})^\top \Gamma^{\mathcal{O}} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}}). \end{aligned}$$

With the symmetry of precision matrices ($\Gamma = \Gamma^\top$), we have

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Gamma (\mathbf{x} - \boldsymbol{\mu}) &= -\frac{1}{2}\mathbf{x}^\top \Gamma (\mathbf{x} - \boldsymbol{\mu}) + \frac{1}{2}\boldsymbol{\mu}^\top \Gamma (\mathbf{x} - \boldsymbol{\mu}) \\ &= -\frac{1}{2}\mathbf{x}^\top \Gamma \mathbf{x} + \frac{1}{2}\mathbf{x}^\top \Gamma \boldsymbol{\mu} + \frac{1}{2}\boldsymbol{\mu}^\top \Gamma \mathbf{x} - \frac{1}{2}\boldsymbol{\mu}^\top \Gamma \boldsymbol{\mu} \\ &= -\frac{1}{2}\mathbf{x}^\top \Gamma \mathbf{x} + \mathbf{x}^\top \Gamma \boldsymbol{\mu} - \frac{1}{2}\boldsymbol{\mu}^\top \Gamma \boldsymbol{\mu}. \end{aligned}$$

Conditional probability - Proof

$$f(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) - \frac{D}{2} \log(2\pi)$$

By using the linear algebra relations

$$-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Gamma} (\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2} \mathbf{x}^\top \boldsymbol{\Gamma} \mathbf{x} + \mathbf{x}^\top \boldsymbol{\Gamma} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Gamma} \boldsymbol{\mu}$$

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} = \frac{\partial}{\partial \mathbf{x}} \mathbf{A}^\top \mathbf{x} = \mathbf{A}, \quad \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x},$$

and by exploiting the derivation chain rule and the symmetry of covariances, we obtain

$$\begin{aligned} -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Gamma} (\mathbf{x} - \boldsymbol{\mu}) &= -\frac{1}{2} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}})^\top \boldsymbol{\Gamma}^{\mathcal{I}} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}}) - \frac{1}{2} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}})^\top \boldsymbol{\Gamma}^{\mathcal{O}} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}}) \\ &\quad - \frac{1}{2} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}})^\top \boldsymbol{\Gamma}^{\mathcal{O}\mathcal{I}} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}}) - \frac{1}{2} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}})^\top \boldsymbol{\Gamma}^{\mathcal{I}\mathcal{O}} (\mathbf{x}^{\mathcal{O}} - \boldsymbol{\mu}^{\mathcal{O}}) \end{aligned}$$

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{x}^{\mathcal{O}}} &= -\boldsymbol{\Gamma}^{\mathcal{O}} \boldsymbol{\mu}^{\mathcal{O}} + \boldsymbol{\Gamma}^{\mathcal{O}\mathcal{I}} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}}) + \boldsymbol{\Gamma}^{\mathcal{O}} \mathbf{x}^{\mathcal{O}} = 0 \\ \iff \hat{\mathbf{x}}^{\mathcal{O}} &= \boldsymbol{\mu}^{\mathcal{O}} - \boldsymbol{\Gamma}^{\mathcal{O}-1} \boldsymbol{\Gamma}^{\mathcal{O}\mathcal{I}} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}^{\mathcal{I}}). \end{aligned}$$

$$\begin{aligned} \boldsymbol{\Sigma}^{-1} &= \begin{bmatrix} \boldsymbol{\Sigma}^{\mathcal{I}}^{-1} + \boldsymbol{\Sigma}^{\mathcal{I}}^{-1} \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} \mathbf{S}^{-1} \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} \boldsymbol{\Sigma}^{\mathcal{I}}^{-1} & -\boldsymbol{\Sigma}^{\mathcal{I}}^{-1} \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} \mathbf{S}^{-1} \\ -\mathbf{S}^{-1} \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} \boldsymbol{\Sigma}^{\mathcal{I}}^{-1} & \mathbf{S}^{-1} \end{bmatrix} \\ \mathbf{S} &= \boldsymbol{\Sigma}^{\mathcal{O}} - \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} \boldsymbol{\Sigma}^{\mathcal{I}}^{-1} \boldsymbol{\Sigma}^{\mathcal{I}\mathcal{O}} \end{aligned}$$

By using the Schur decomposition, we can see that

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

Conditional probability - Proof

$$\Sigma^{-1} = \begin{bmatrix} \Sigma^{\mathcal{I}-1} + \Sigma^{\mathcal{I}-1} \Sigma^{\mathcal{IO}} S^{-1} \Sigma^{\mathcal{OI}} \Sigma^{\mathcal{I}-1} & -\Sigma^{\mathcal{I}-1} \Sigma^{\mathcal{IO}} S^{-1} \\ -S^{-1} \Sigma^{\mathcal{OI}} \Sigma^{\mathcal{I}-1} & S^{-1} \end{bmatrix}$$
$$S = \Sigma^{\mathcal{O}} - \Sigma^{\mathcal{OI}} \Sigma^{\mathcal{I}-1} \Sigma^{\mathcal{IO}}$$

The associated covariance matrix $\hat{\Sigma}^{\mathcal{O}}$ measuring the error of this estimate is given by the inverse of the Hessian matrix \mathbf{H} . We have

$$\mathbf{H} = \frac{\partial^2 f}{\partial \mathbf{x}^{\mathcal{O}} \mathbf{x}^{\mathcal{O}\top}} = \boldsymbol{\Gamma}^{\mathcal{O}} \quad \Rightarrow \quad \hat{\Sigma}^{\mathcal{O}} = \boldsymbol{\Gamma}^{\mathcal{O}-1}.$$

We can then see that

$$\hat{\Sigma}^{\mathcal{O}} = S = \Sigma^{\mathcal{O}} - \Sigma^{\mathcal{OI}} \Sigma^{\mathcal{I}-1} \Sigma^{\mathcal{IO}}.$$

Note that in many cases, evaluating the conditional distribution with precision matrices is computationally more efficient than with covariance matrices.

Conditional probability - Summary

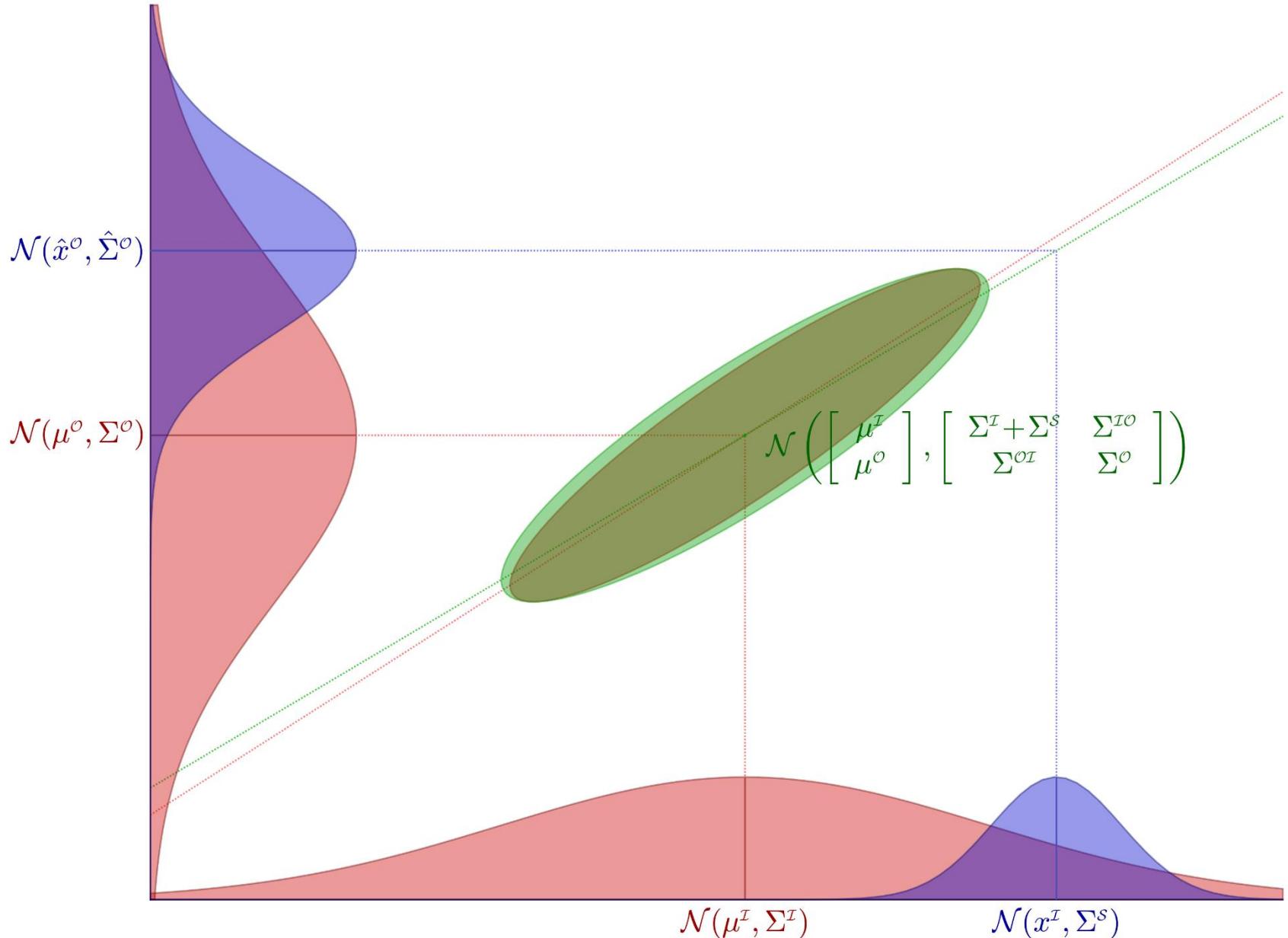
If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we have that $\mathbf{x}^{\mathcal{O}} | \mathbf{x}^{\mathcal{I}} \sim \mathcal{N}(\hat{\mathbf{x}}^{\mathcal{O}}, \hat{\boldsymbol{\Sigma}}^{\mathcal{O}})$, with parameters

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

If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Gamma}^{-1})$, we have that $\mathbf{x}^{\mathcal{O}} | \mathbf{x}^{\mathcal{I}} \sim \mathcal{N}(\hat{\mathbf{x}}^{\mathcal{O}}, \hat{\boldsymbol{\Gamma}}^{\mathcal{O}-1})$, with parameters

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

Gaussian conditioning with uncertain inputs



Gaussian conditioning with uncertain inputs

For inputs corrupted by noise $\boldsymbol{\epsilon}^{\mathcal{I}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{\textcolor{red}{s}})$, we have that $\mathbf{x}^{\mathcal{O}} | \mathbf{x}^{\mathcal{I}} + \boldsymbol{\epsilon}^{\mathcal{I}} \sim \mathcal{N}(\hat{\mathbf{x}}^{\mathcal{O}}, \hat{\boldsymbol{\Sigma}}^{\mathcal{O}})$ with parameters

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

This can for example be shown by redefining the joint distribution $\boldsymbol{\Sigma}$ with

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}^{\mathcal{I}} + \boldsymbol{\Sigma}^{\textcolor{red}{s}} & \boldsymbol{\Sigma}^{\mathcal{I}\mathcal{O}} \\ \boldsymbol{\Sigma}^{\mathcal{O}\mathcal{I}} & \boldsymbol{\Sigma}^{\mathcal{O}} \end{bmatrix},$$

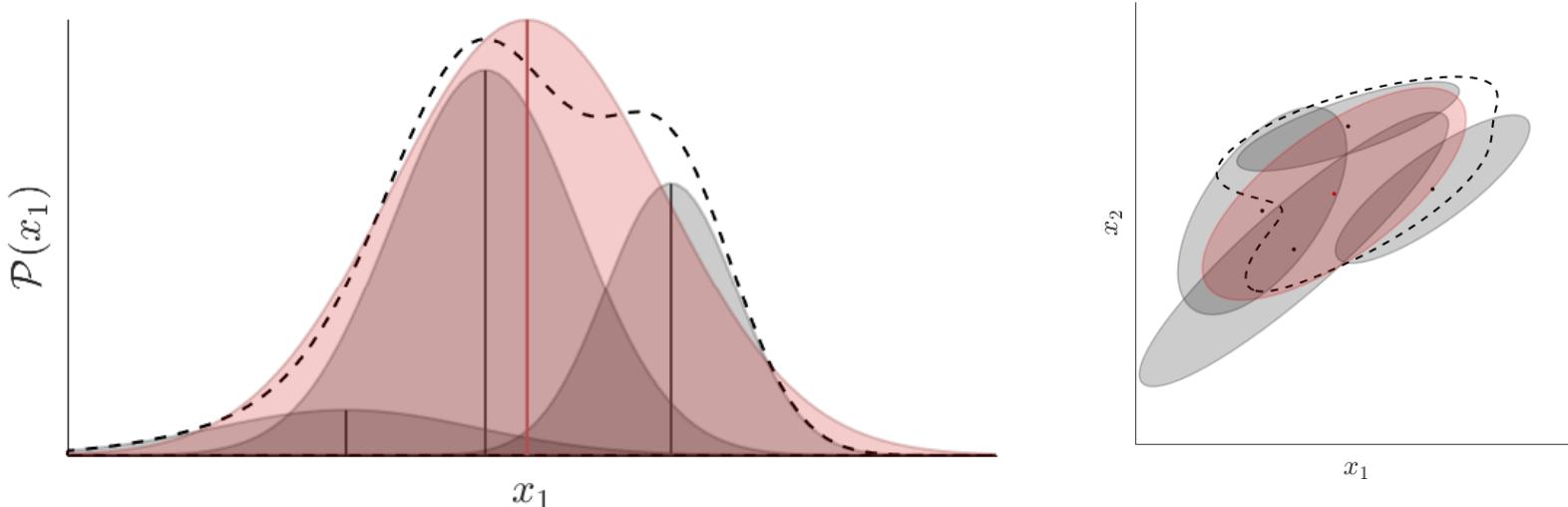
and applying the same formulas as for standard Gaussian conditioning.

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 estimate of a mixture - Proof

The result can be demonstrated by developing the expressions

$$\begin{aligned} \mathbb{E}(\mathbf{x}) &= \boldsymbol{\mu}, \quad \boldsymbol{\Sigma} = \text{cov}(\mathbf{x}) = \mathbb{E}(\mathbf{x}\mathbf{x}^\top) - \mathbb{E}(\mathbf{x})\mathbb{E}(\mathbf{x}^\top) = \mathbb{E}(\mathbf{x}\mathbf{x}^\top) - \boldsymbol{\mu}\boldsymbol{\mu}^\top \\ \iff \mathbb{E}(\mathbf{x}\mathbf{x}^\top) &= \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^\top. \end{aligned}$$

By considering datapoints \mathbf{x} distributed with a mixture of Gaussians

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

the mean is computed as

$$\begin{aligned} \boldsymbol{\mu} &= \mathbb{E}(\mathbf{x}) = \int \mathbf{x} \mathcal{P}(\mathbf{x}) d\mathbf{x} = \int \mathbf{x} \sum_{i=1}^K h_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x} \\ \mathbb{E}(\mathbf{x}) &= \int \mathbf{x} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x} \\ &\quad = \sum_{i=1}^K h_i \int \mathbf{x} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x} \\ &\quad = \sum_{i=1}^K h_i \boldsymbol{\mu}_i. \end{aligned}$$

Gaussian estimate of a mixture of Gaussians

By noting that

$$\mathbb{E}(\mathbf{x}\mathbf{x}^\top) = \Sigma + \mu\mu^\top$$

$$\begin{aligned}\mathbb{E}(\mathbf{x}\mathbf{x}^\top) &= \int \mathbf{x}\mathbf{x}^\top \mathcal{P}(\mathbf{x}) d\mathbf{x} \\ &= \int \sum_{i=1}^K h_i \mathbf{x}\mathbf{x}^\top \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x} \\ &= \sum_{i=1}^K h_i \int \mathbf{x}\mathbf{x}^\top \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x} \\ \mathbb{E}(\mathbf{x}\mathbf{x}^\top) &= \int \mathbf{x}\mathbf{x}^\top \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x} \\ &\stackrel{\text{red}}{=} \sum_{i=1}^K h_i \left(\boldsymbol{\Sigma}_i + \boldsymbol{\mu}_i \boldsymbol{\mu}_i^\top \right),\end{aligned}$$

the covariance is then computed as

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

Locally weighted regression (LWR)

**Matlab codes: demo_LWR01.m,
demo_DMP01.m**

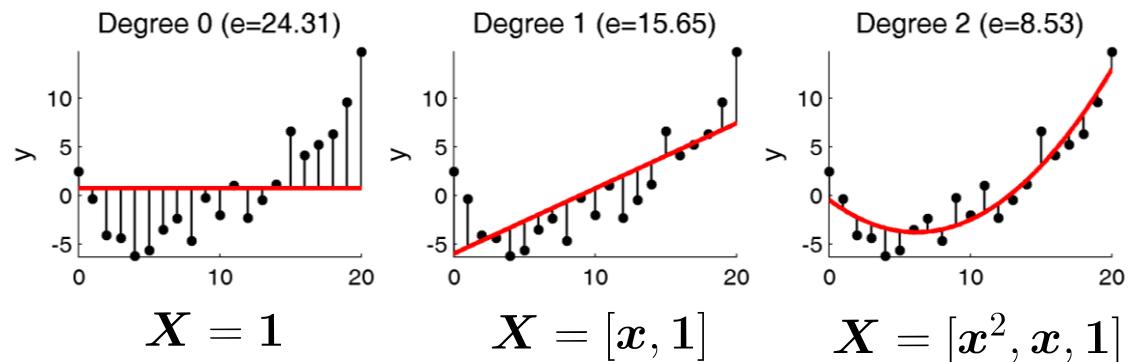
[C.G. Atkeson, A.W. Moore, and S. Schaal. Locally weighted learning for control. Artificial Intelligence Review, 11(1-5):75–113, 1997]

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

Previous lecture on linear regression

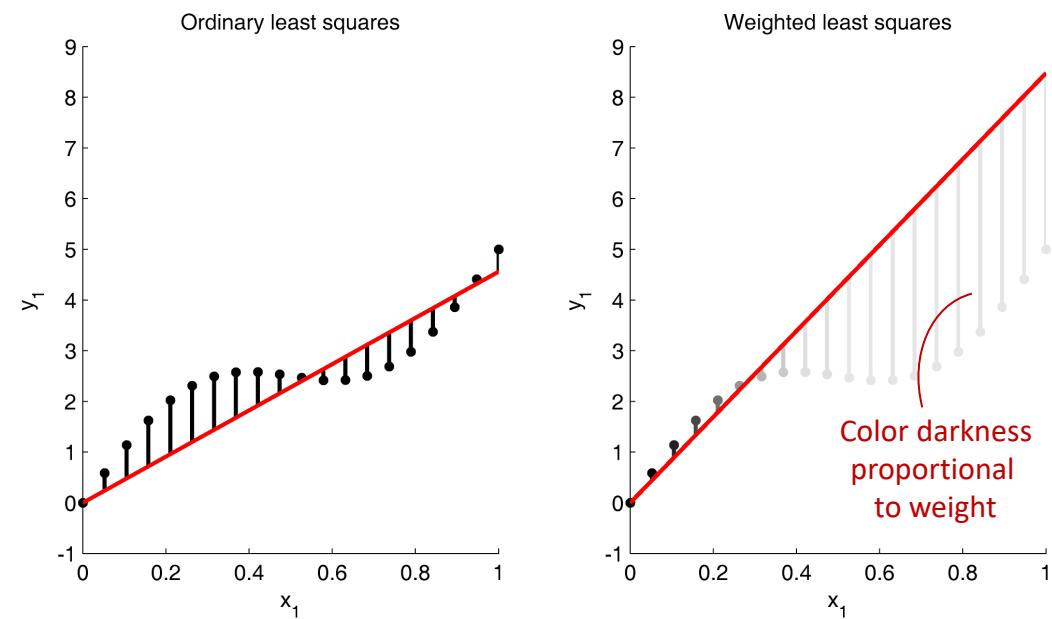
$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} (\mathbf{Y} - \mathbf{X}\mathbf{A})^\top (\mathbf{Y} - \mathbf{X}\mathbf{A})$$

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



$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} (\mathbf{Y} - \mathbf{X}\mathbf{A})^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\mathbf{A})$$

$$= (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$



Locally weighted regression (LWR)

Locally weighted regression (LWR) is a direct extension of the weighted least squares formulation in which K weighted regressions are performed on the same dataset $\{\mathbf{X}^I, \mathbf{X}^O\}$.

It aims at splitting a nonlinear problem so that it can be solved locally by linear regression.

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

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

or in its rescaled form as

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

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

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

$$\begin{aligned}\mathbf{X}^{\mathcal{I}} &= [t_1, t_2, \dots, t_N]^{\top} \\ \hat{\mathbf{A}}_k &= (\mathbf{X}^{\mathcal{I}\top} \mathbf{W}_k \mathbf{X}^{\mathcal{I}})^{-1} \mathbf{X}^{\mathcal{I}\top} \mathbf{W}_k \mathbf{X}^{\mathcal{O}}\end{aligned}$$

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

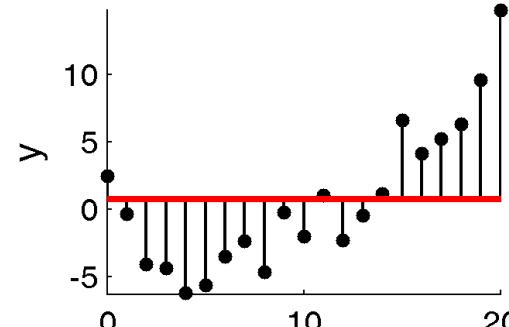
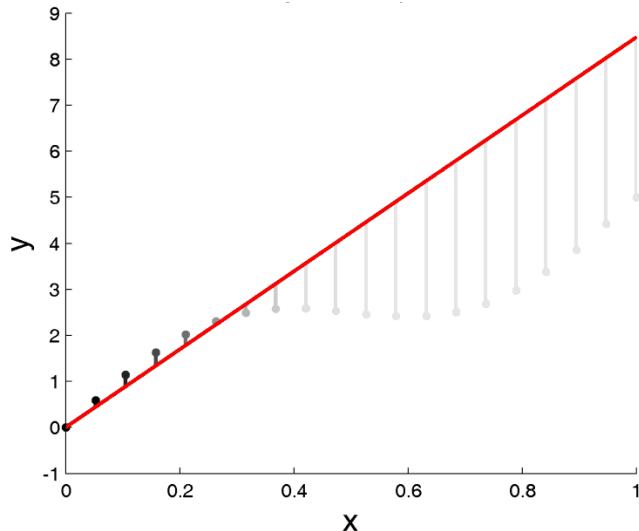
can then be used to evaluate $\hat{\mathbf{A}}_k$. The result can finally be used to compute

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

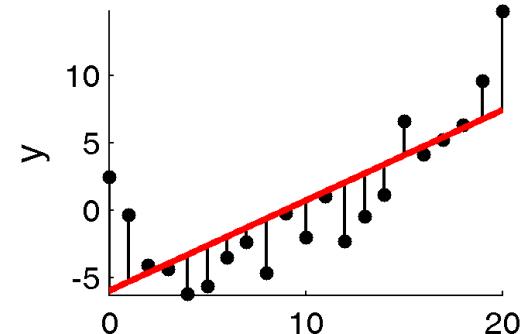
Multiple variants exist, including online estimation within a recursive formulation [Schaal'98], Bayesian treatments of LWR [Ting'08], and also extensions such as **locally weighted projection regression** (LWPR) that exploit partial least squares to handle redundant and/or irrelevant inputs, with an online algorithm to estimate the model parameters [Vijayakumar'05].

Locally weighted regression (LWR)

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

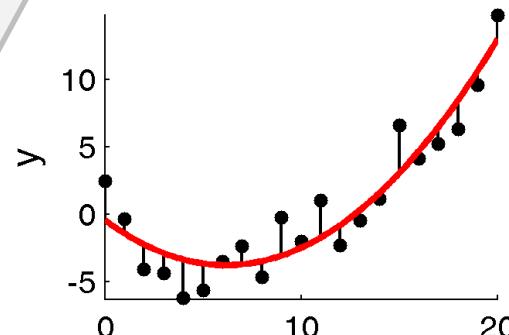


$$\mathbf{X} = 1$$



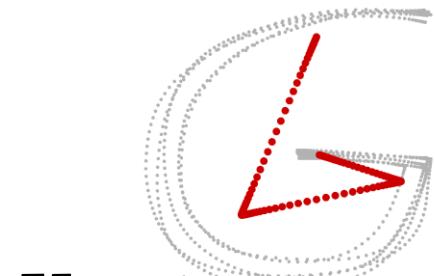
$$\mathbf{X} = [x, 1]$$

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

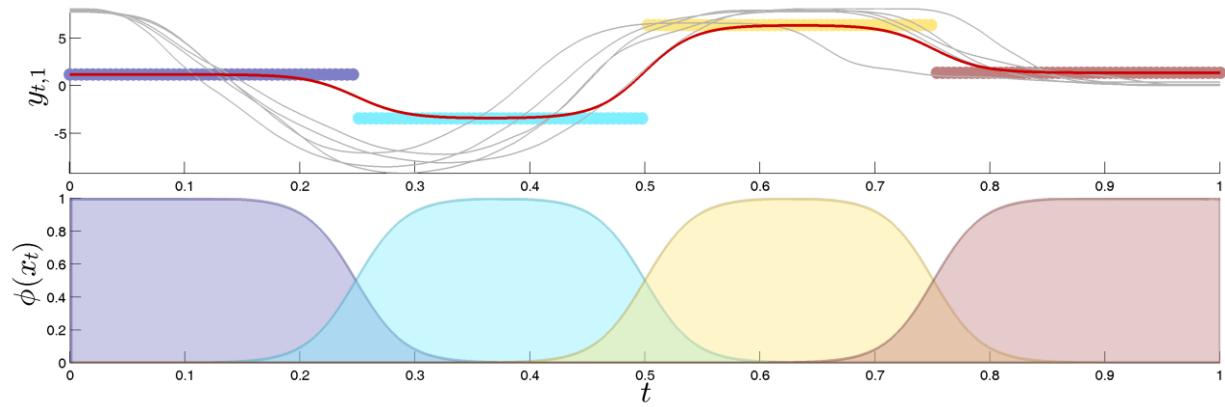


$$\mathbf{X} = [x^2, x, 1]$$

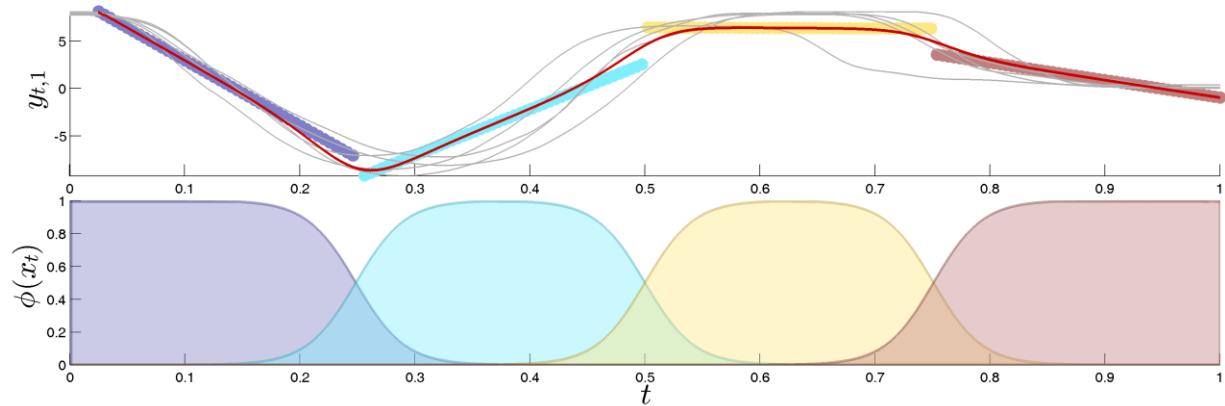
Locally weighted regression (LWR)



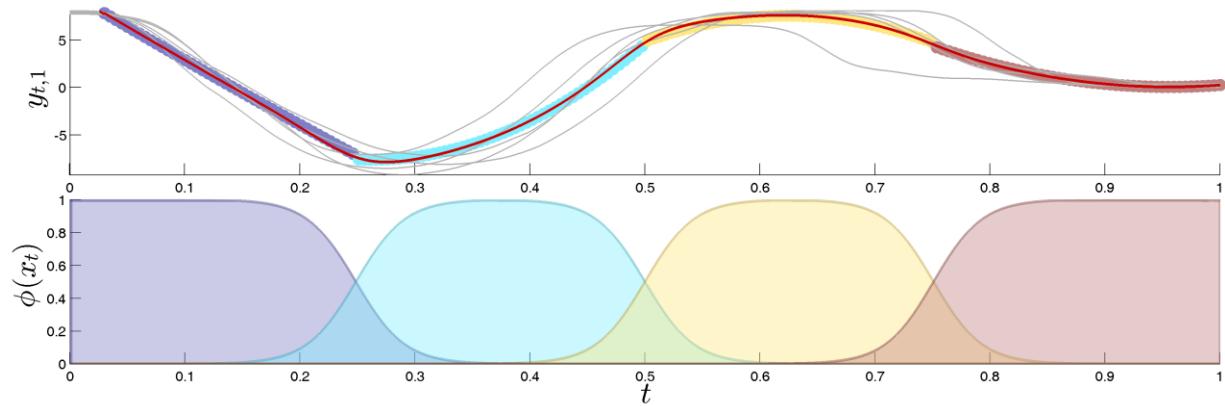
$$X = 1$$



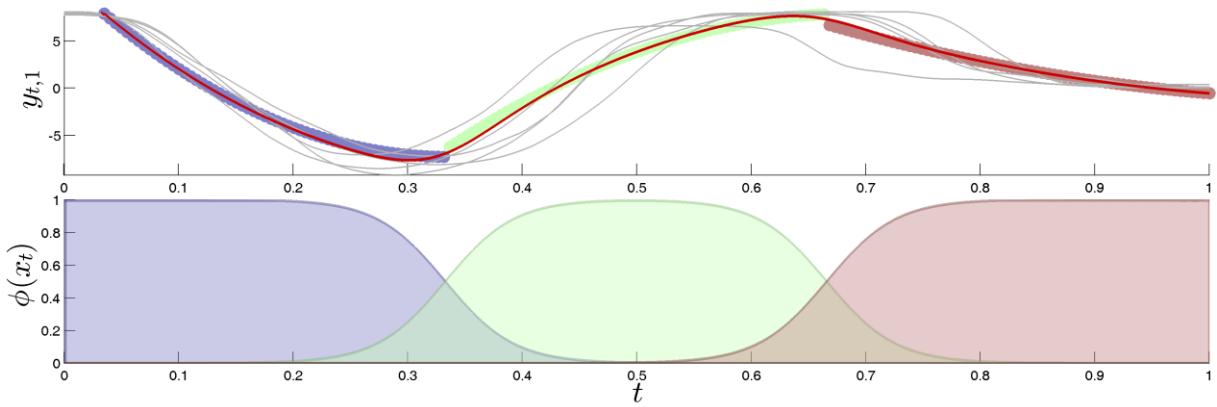
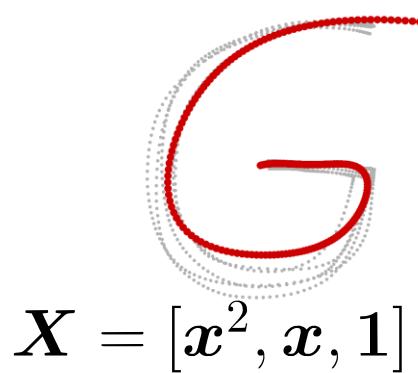
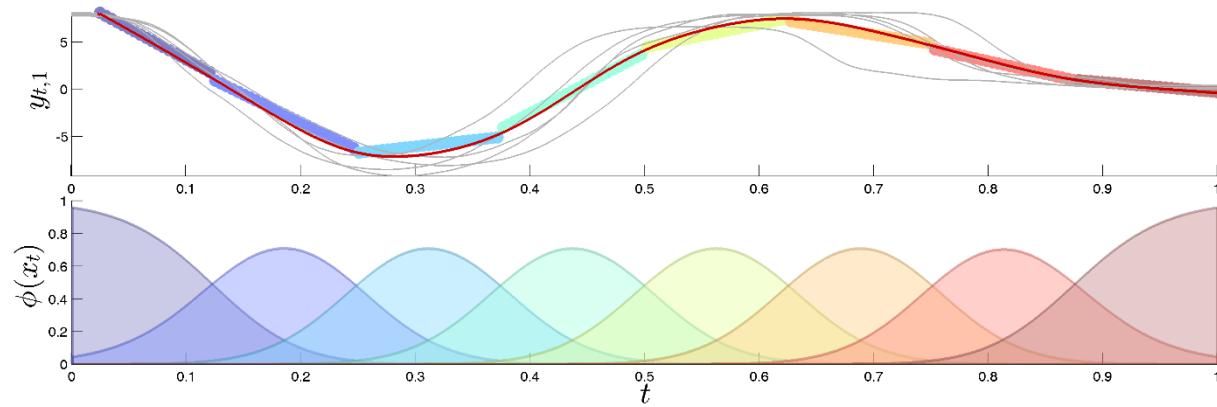
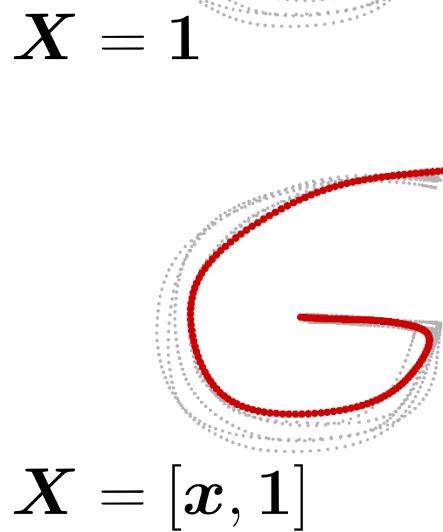
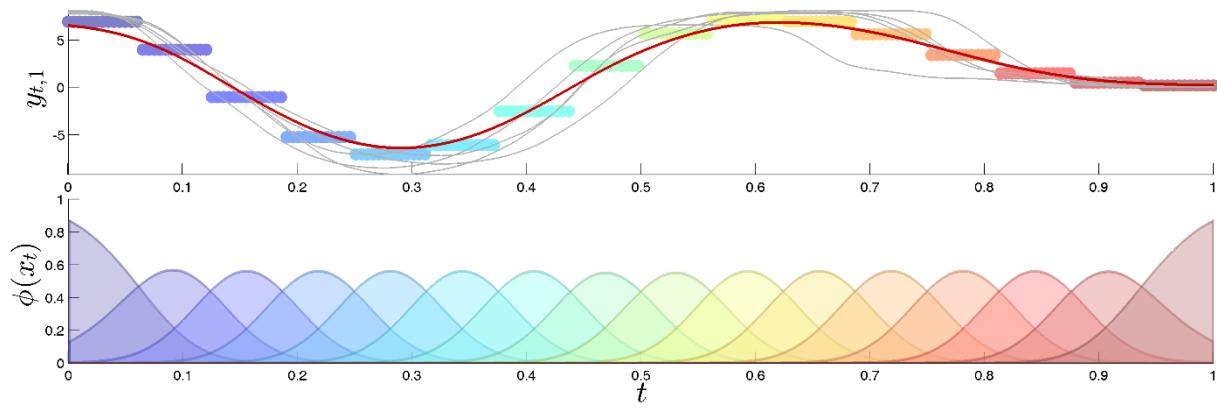
$$X = [x, 1]$$



$$X = [x^2, x, 1]$$

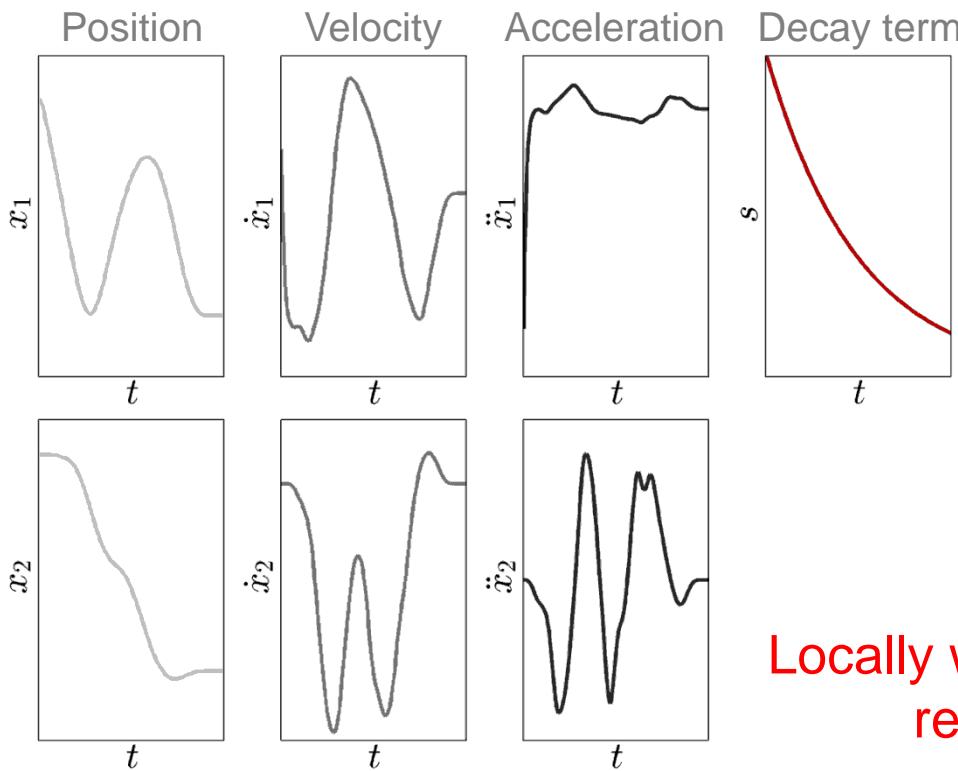


Locally weighted regression (LWR)

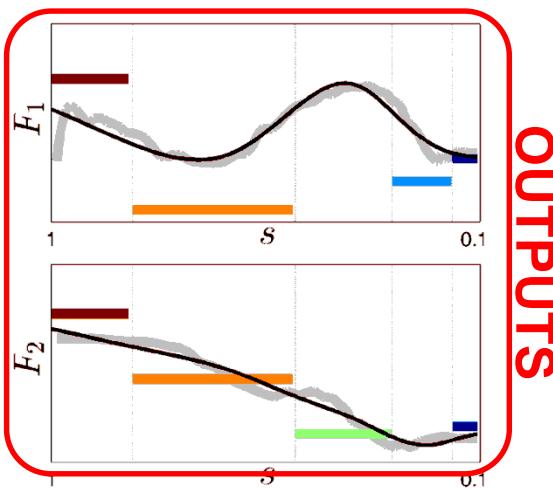
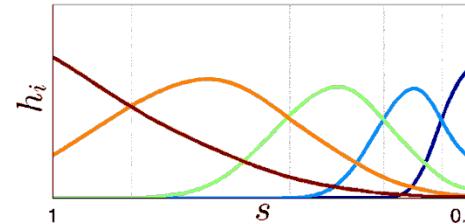


Dynamical movement primitives (DMP)

INPUTS



Set of basis functions



Locally weighted
regression
(LWR)

$$\tau \ddot{x} = \kappa^p [x_T - x] - \kappa^\nu \dot{x} + f(s), \quad f(s) = s \sum_{i=1}^K h_i(s) f_i$$

$$\tau \dot{s} = -\alpha s$$

[Ijspeert, Nakanishi and Schaal, NIPS'2003]

[Ijspeert, Nakanishi, Pastor, Hoffmann and Schaal, Neural Computation 25(2), 2013]

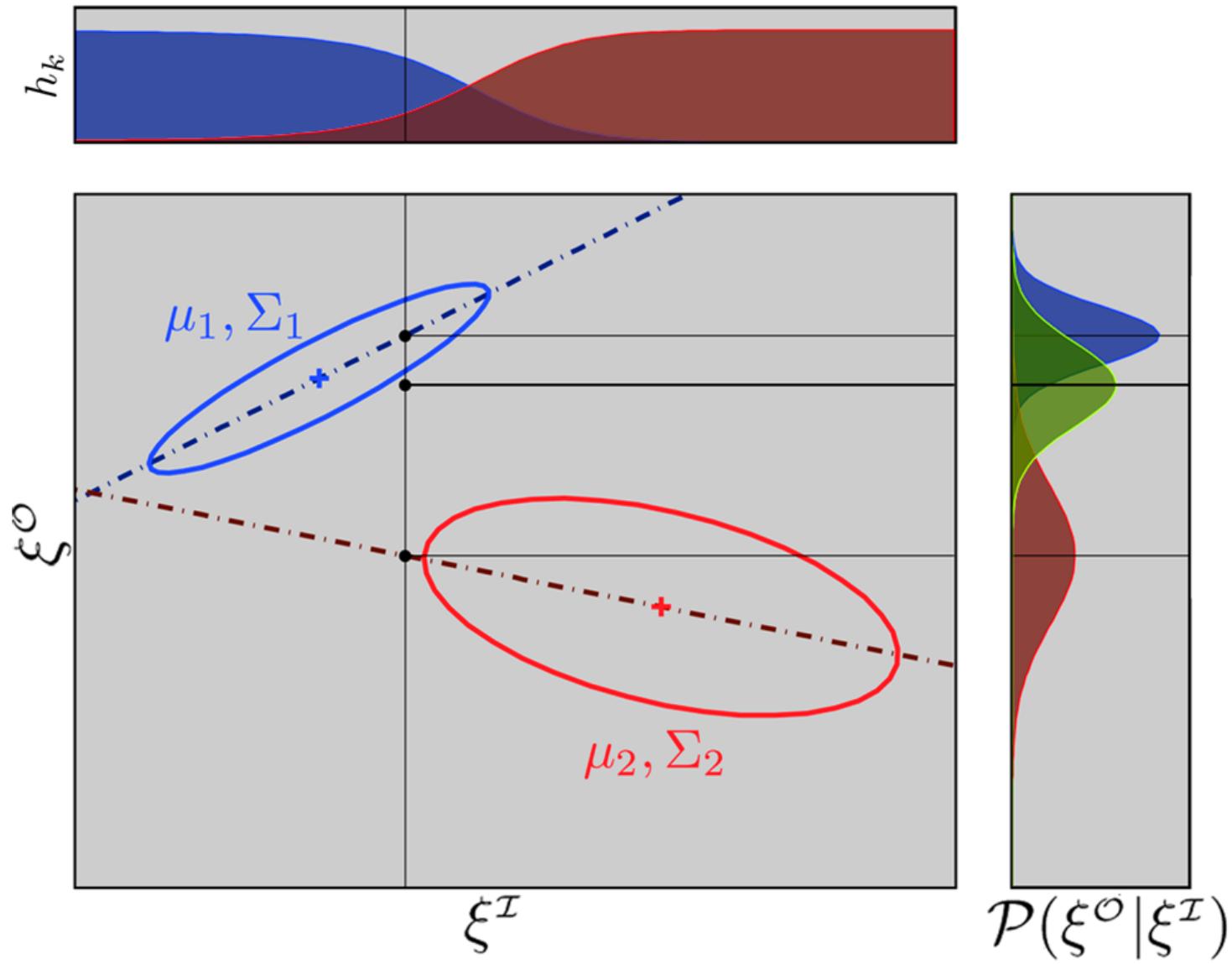
Gaussian mixture regression (GMR)

Matlab codes: `demo_GMR01.m`
`demo_GMR_polyFit01.m,`
`demo_DMP_GMR01.m`

[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 S. A tutorial on task-parameterized movement learning and retrieval. Intelligent Service Robotics 9(1):1–29, 2016]

Gaussian mixture regression (GMR)



Gaussian mixture regression (GMR)

- Gaussian mixture regression (GMR) is a nonlinear regression technique that does not model the regression function directly, but instead first models the **joint probability density of input-output data** in the form of a Gaussian mixture model (GMM).
- The computation relies on **linear transformation and conditioning properties** of multivariate normal distributions.
- GMR provides a regression approach in which **multivariate output distributions can be computed in an online manner**, with a computation time **independent of the number of datapoints** used to train the model, by exploiting the learned joint density model.
- In GMR, **both input and output variables can be multivariate**, and after learning, **any subset of input-output dimensions can be selected** for regression. This can for example be exploited to handle different sources of missing data, where expectations on the remaining dimensions can be computed as a multivariate distribution.

Gaussian mixture regression (GMR)

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

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

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

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

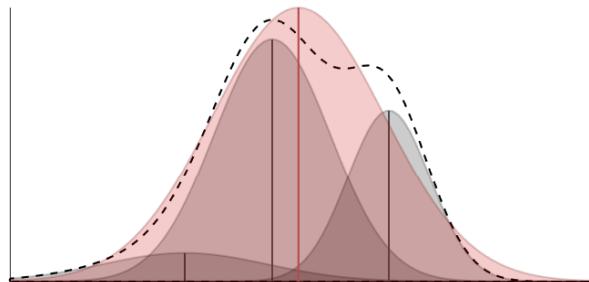
$$\hat{\boldsymbol{\Sigma}}_i^{\mathcal{O}} = \boldsymbol{\Sigma}_i^{\mathcal{O}} - \boldsymbol{\Sigma}_i^{\mathcal{O}\mathcal{I}} \boldsymbol{\Sigma}_i^{\mathcal{I}}^{-1} \boldsymbol{\Sigma}_i^{\mathcal{I}\mathcal{O}}$$

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

computed with the marginal

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

Gaussian mixture regression (GMR)



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

$$\begin{aligned}\hat{\boldsymbol{\mu}}_i^{\mathcal{O}} &= \boldsymbol{\mu}_i^{\mathcal{O}} + \boldsymbol{\Sigma}_i^{\mathcal{O}\mathcal{I}} \boldsymbol{\Sigma}_i^{\mathcal{I}-1} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}_i^{\mathcal{I}}) \\ \hat{\boldsymbol{\Sigma}}_i^{\mathcal{O}} &= \boldsymbol{\Sigma}_i^{\mathcal{O}} - \boldsymbol{\Sigma}_i^{\mathcal{O}\mathcal{I}} \boldsymbol{\Sigma}_i^{\mathcal{I}-1} \boldsymbol{\Sigma}_i^{\mathcal{I}\mathcal{O}}\end{aligned}$$

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

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

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

Gaussian mixture regression (GMR) - Proof

The above can be demonstrated by computing

$$\hat{\boldsymbol{\mu}}^{\mathcal{O}} = \mathbb{E}(\boldsymbol{x}^{\mathcal{O}} | \boldsymbol{x}^{\mathcal{I}}),$$

$$\hat{\boldsymbol{\Sigma}}^{\mathcal{O}} = \text{cov}(\boldsymbol{x}^{\mathcal{O}} | \boldsymbol{x}^{\mathcal{I}}) = \mathbb{E}(\boldsymbol{x}^{\mathcal{O}} \boldsymbol{x}^{\mathcal{O}\top} | \boldsymbol{x}^{\mathcal{I}}) - \mathbb{E}(\boldsymbol{x}^{\mathcal{O}} | \boldsymbol{x}^{\mathcal{I}}) \mathbb{E}(\boldsymbol{x}^{\mathcal{O}\top} | \boldsymbol{x}^{\mathcal{I}}).$$

The conditional mean can be computed as

$$\hat{\boldsymbol{\mu}}^{\mathcal{O}} = \mathbb{E}(\boldsymbol{x}^{\mathcal{O}} | \boldsymbol{x}^{\mathcal{I}}) = \sum_{i=1}^K h_i \hat{\boldsymbol{\mu}}_i^{\mathcal{O}}.$$

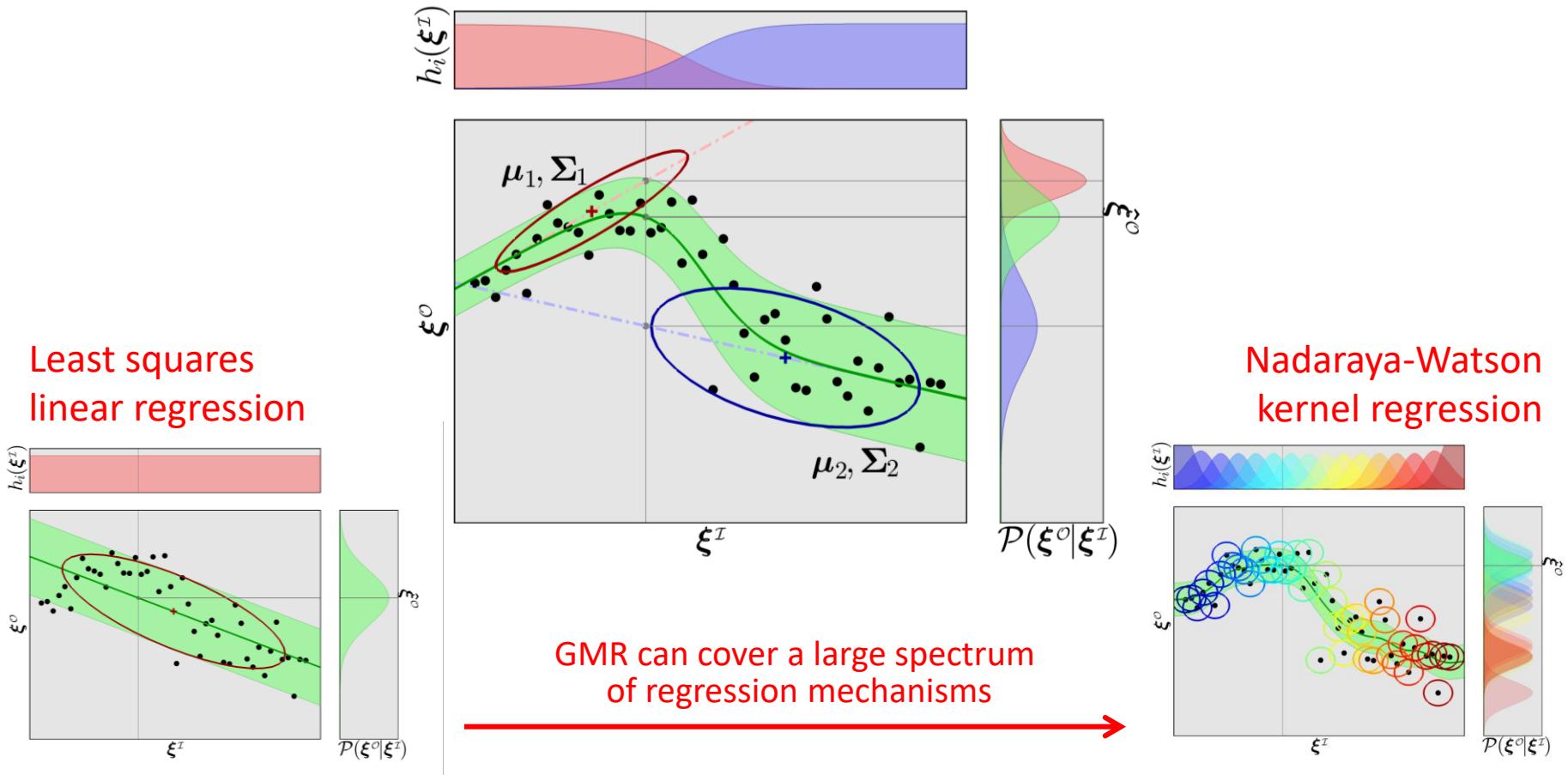
In order to evaluate the covariance, we first note that

$$\mathbb{E}(\boldsymbol{x}^{\mathcal{O}} \boldsymbol{x}^{\mathcal{O}\top} | \boldsymbol{x}^{\mathcal{I}}) = \sum_{i=1}^K h_i \hat{\boldsymbol{\Sigma}}_i^{\mathcal{O}} + \sum_{i=1}^K h_i \hat{\boldsymbol{\mu}}_i^{\mathcal{O}} \hat{\boldsymbol{\mu}}_i^{\mathcal{O}\top}.$$

We then have

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

Gaussian mixture regression (GMR)



Both $\xi^{\mathcal{I}}$ and $\xi^{\mathcal{O}}$ can be multidimensional

$P(\xi^{\mathcal{I}}, \xi^{\mathcal{O}})$ encoded in **Gaussian mixture model (GMM)**

$P(\xi^{\mathcal{O}}|\xi^{\mathcal{I}})$ retrieved by **Gaussian mixture regression (GMR)**

GMR with uncertain inputs

For inputs with noise $\boldsymbol{\epsilon}^{\mathcal{I}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{\mathcal{S}})$, we have that $\mathbf{x}^{\mathcal{O}} | \mathbf{x}^{\mathcal{I}} + \boldsymbol{\epsilon}^{\mathcal{I}} \sim \mathcal{N}(\hat{\boldsymbol{\mu}}^{\mathcal{O}}, \hat{\boldsymbol{\Sigma}}^{\mathcal{O}})$, with parameters

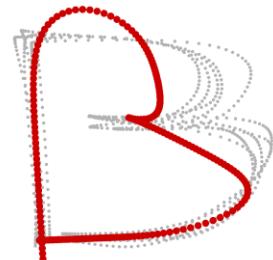
$$\begin{aligned}\hat{\boldsymbol{\mu}}^{\mathcal{O}} &= \sum_{i=1}^K h_i \hat{\boldsymbol{\mu}}_i^{\mathcal{O}}, \\ \hat{\boldsymbol{\Sigma}}^{\mathcal{O}} &= \sum_{i=1}^K h_i \left(\hat{\boldsymbol{\Sigma}}_i^{\mathcal{O}} + \hat{\boldsymbol{\mu}}_i^{\mathcal{O}} \hat{\boldsymbol{\mu}}_i^{\mathcal{O}\top} \right) - \hat{\boldsymbol{\mu}}^{\mathcal{O}} \hat{\boldsymbol{\mu}}^{\mathcal{O}\top},\end{aligned}$$

computed with

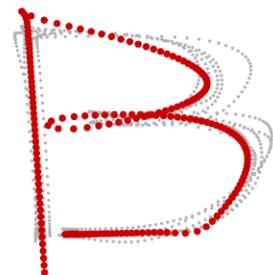
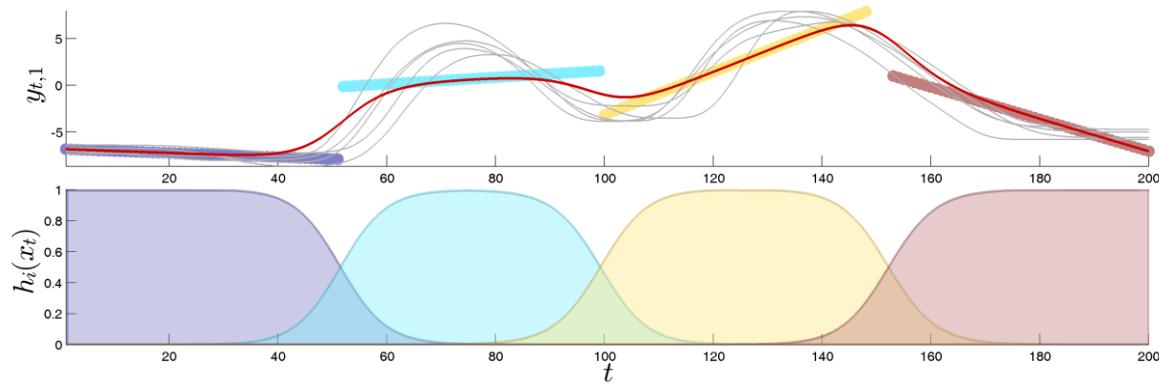
$$\begin{aligned}\hat{\boldsymbol{\mu}}_i^{\mathcal{O}} &= \boldsymbol{\mu}_i^{\mathcal{O}} + \boldsymbol{\Sigma}_i^{\mathcal{O}\mathcal{I}} (\boldsymbol{\Sigma}_i^{\mathcal{I}} + \boldsymbol{\Sigma}^{\mathcal{S}})^{-1} (\mathbf{x}^{\mathcal{I}} - \boldsymbol{\mu}_i^{\mathcal{I}}), \\ \hat{\boldsymbol{\Sigma}}_i^{\mathcal{O}}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

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

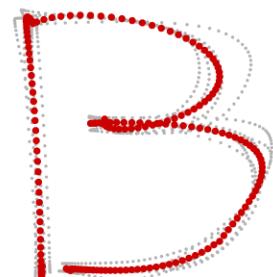
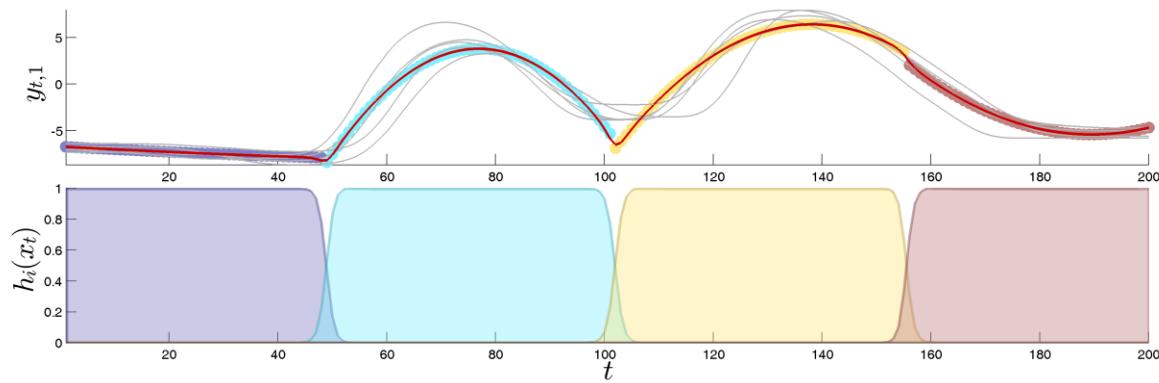
GMR for smooth piecewise polynomial fitting



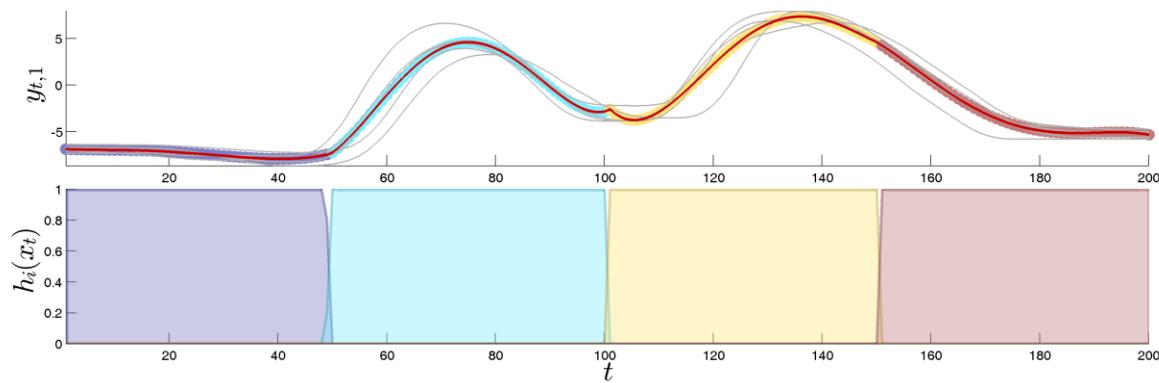
$$X = x$$



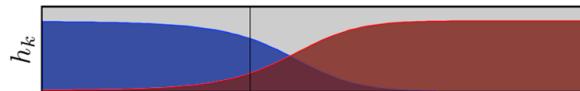
$$X = [x^2, x]$$



$$X = [x^3, x^2, x]$$

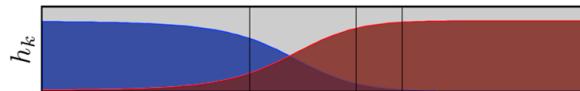
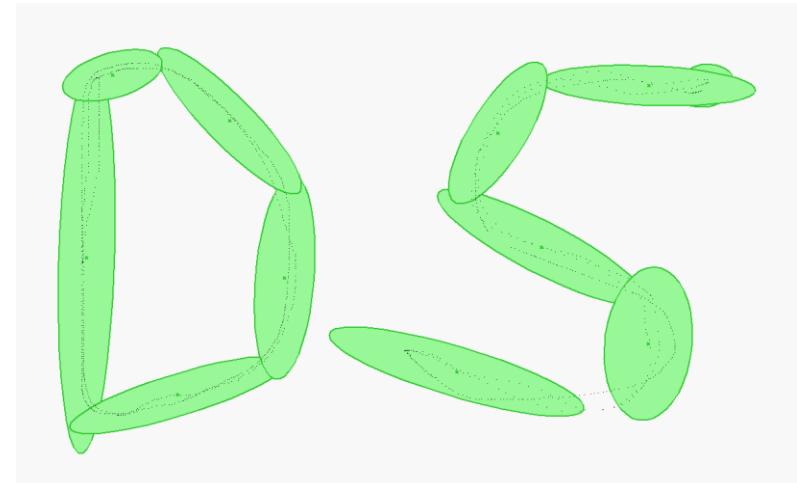
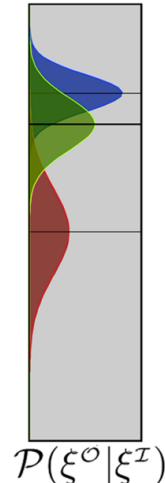
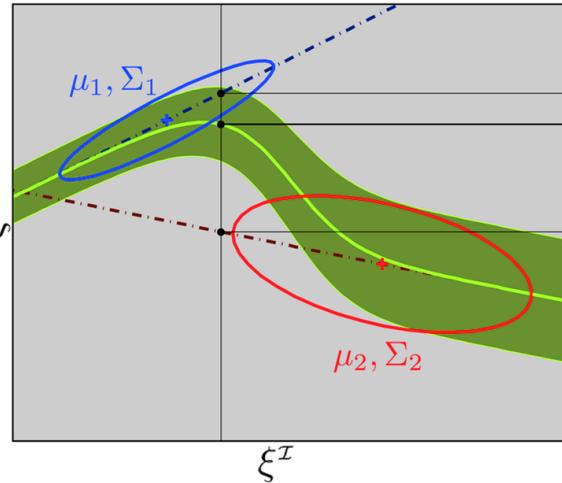


Gaussian mixture regression - Examples



$$\xi^{\mathcal{I}} = \mathbf{t}, \quad \xi^{\mathcal{O}} = \mathbf{x}$$

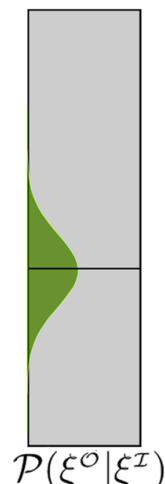
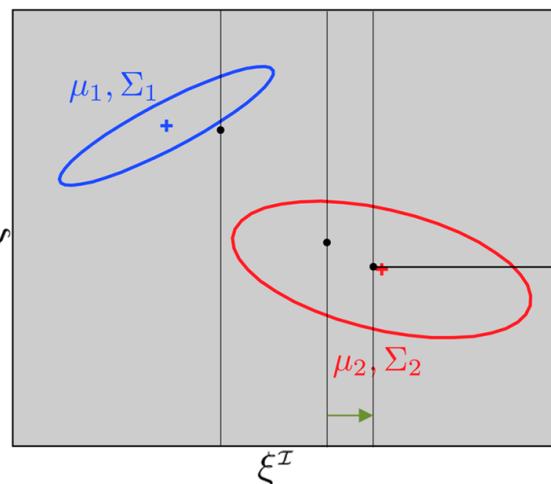
[Calinon, Guenter and Billard,
IEEE Trans. on SMC-B 37(2), 2007]



$$\xi^{\mathcal{I}} = \mathbf{x}, \quad \xi^{\mathcal{O}} = \dot{\mathbf{x}}$$

With expectation-maximization (EM):
(maximizing log-likelihood)

[Hersch, Guenter, Calinon and Billard,
IEEE Trans. on Robotics 24(6), 2008]

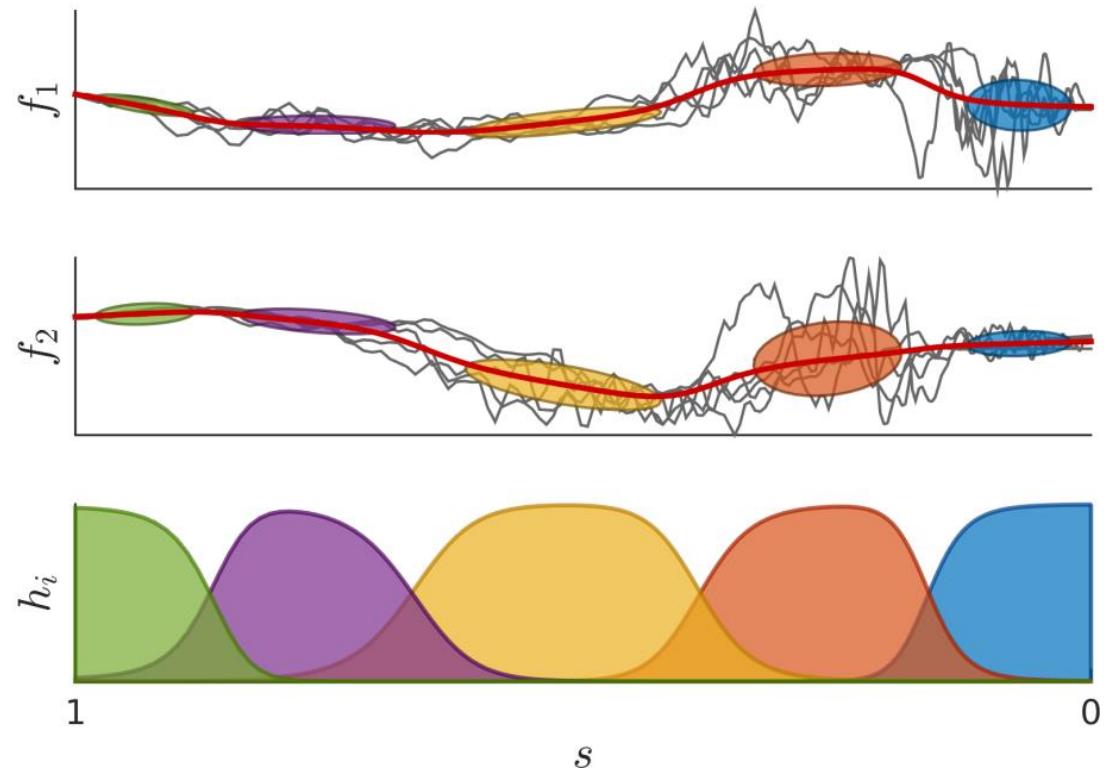
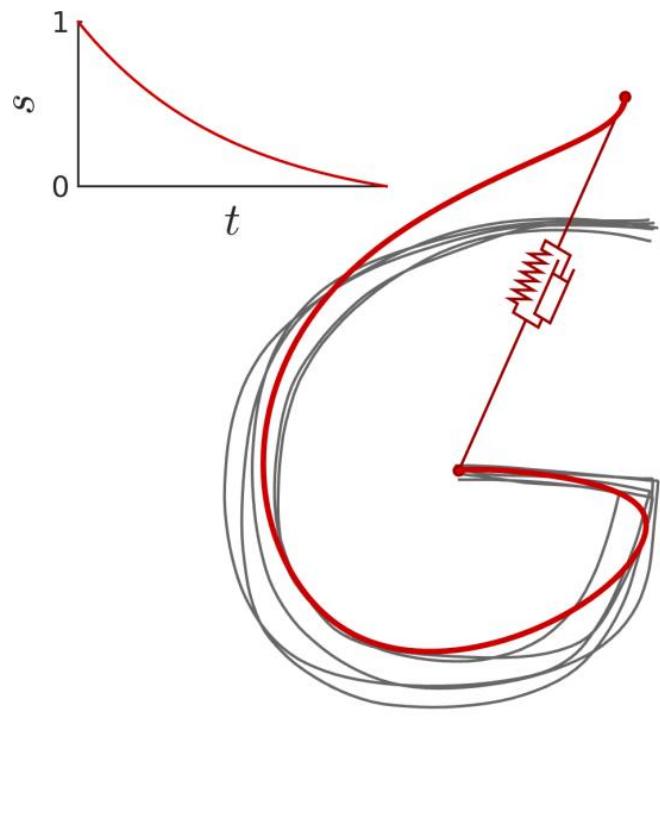


With quadratic programming solver:
(maximizing log-likelihood s.t. stability constraints)

[Khansari-Zadeh and Billard,
IEEE Trans. on Robotics 27(5), 2011]

Dynamical movement primitives with GMR

Learning of $\mathcal{P}(s, \mathbf{x})$ and retrieval of $\mathcal{P}(\mathbf{x}|s)$



Main references

Regression

F. Stulp and O. Sigaud. Many regression algorithms, one unified model – a review. *Neural Networks*, 69:60–79, 2015

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

Recursive formulation of LWR

S. Schaal and C.G. Atkeson. Constructive incremental learning from only local information. *Neural Computation* 10(8):2047–2084, 1998.

Bayesian formulation of LWR

J. Ting, M. Kalakrishnan, S. Vijayakumar and S. Schaal. Bayesian kernel shaping for learning control. In: *Advances in Neural Information Processing Systems (NIPS)*, pp 1673–1680, 2008.

Main references

DMP

A. Ijspeert, J. Nakanishi, P. Pastor, H. Hoffmann, and S. Schaal. Dynamical movement primitives: Learning attractor models for motor behaviors. *Neural Computation*, 25(2):328–373, 2013

LWPR

S. Vijayakumar, A. D’souza and S. Schaal. Incremental online learning in high dimensions. *Neural Computation* 17(12):2602–2634, 2005

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 S. A tutorial on task-parameterized movement learning and retrieval. *Intelligent Service Robotics* 9(1):1–29, 2016