

Non-Linear Regression

Gaussian Mixture Regression (GMR)





Gaussian Mixture Regression (GMR): Principle

Given a multimodal dataset

$$x \in \mathbb{R}^{N_x}, y \in \mathbb{R}^{N_y}$$

x,y are random vectors with dimension N_x and N_y respectively

- 1) Estimate the *joint* density, p(x,y)
- 2) Estimate y given x by estimating the conditional density, p(y/x)

We need to decide on a probabilistic model for p(x,y). *GMR assumes that* p(x,y) is modeled by a mixture of Gaussians

$$p(x,y) = \sum_{k=1}^{K} \alpha_k \cdot p(x,y;\mu^k, \Sigma^k), \quad \text{with } p(x,y;\mu^k, \Sigma^k) = N(\mu^k, \Sigma^k)$$

 μ^k , Σ^k : mean and covariance matrix of Gaussian component k.



GMR: with a single Gaussian component (k=1)

 $p(x, y) \sim \mathcal{N}(\mu, \Sigma)$; μ and Σ are maximum likelihood estimates

$$\mu = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}$$
 vector with dimensionality $N_x + N_y$ $\mu_x = \frac{1}{M} \sum_{m=1}^M x^m$, $\mu_y = \frac{1}{M} \sum_{m=1}^M y^m$

$$\Sigma = \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_y \end{bmatrix}$$
 symetric matrix with dimensionality $(N_x + N_y) \times (N_x + N_y)$

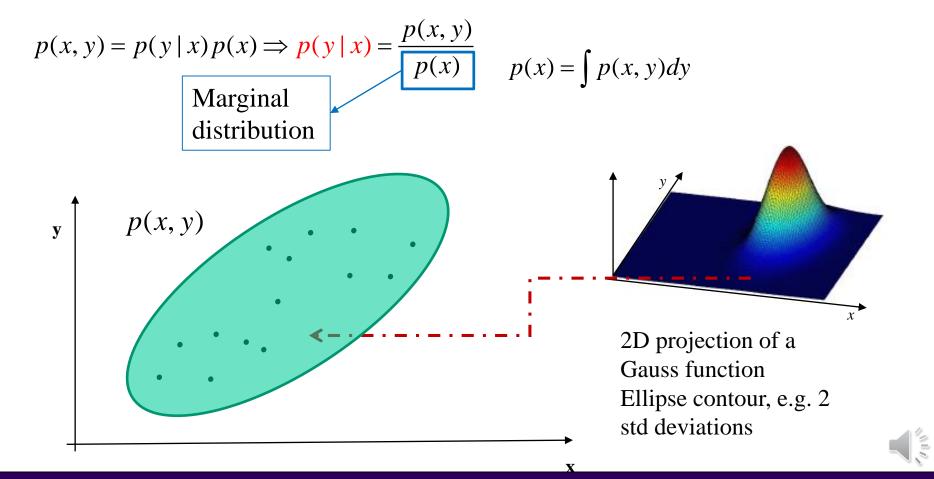
$$\Sigma_{x} = \frac{1}{M} \sum_{m=1}^{M} \left(x^{m} - \mu_{x} \right) \left(x^{m} - \mu_{x} \right)^{T}, \ \Sigma_{xy} = \frac{1}{M} \sum_{m=1}^{M} \left(x^{m} - \mu_{x} \right) \left(y^{m} - \mu_{y} \right)^{T}$$





GMR: with a single Gaussian component (k=1)

- 1) We first estimate the *joint* density, p(x,y), across pairs of datapoints with a single Gauss distribution.
- 2) Then, we compute p(y|x) in order to estimate y for a query point x.



4



GMR: with a single Gaussian component (k=1)

The marginals and the conditionals of a joint Gauss distribution are also Gauss distributions¹

Joint Density

$$p(x, y) \sim \mathcal{N}(x, y \mid \mu, \Sigma)$$

$$\mu = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_y \end{bmatrix} \qquad p(y) \sim \mathcal{N}(y \mid \mu_y, \Sigma_y)$$

Marginals

$$p(x) \sim \mathcal{N}(x \mid \mu_x, \Sigma_x)$$

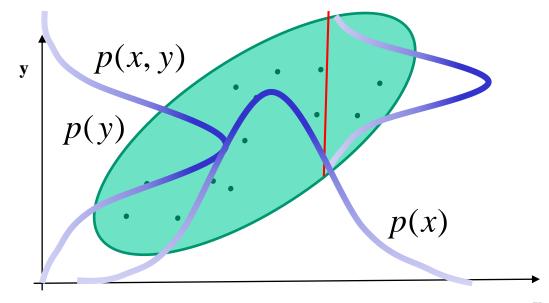
$$p(y) \sim \mathcal{N}(y | \mu_y, \Sigma_y)$$

Conditional

$$p(y \mid x) \sim \mathcal{N}(y; \mu_{y|x}, \Sigma_{y|x})$$

$$\mu_{y|x} = \mu_y + \sum_{yx} \sum_{x}^{-1} (x - \mu_x)$$

$$\Sigma_{y|x} = \Sigma_y - \Sigma_{yx} \Sigma_x^{-1} \Sigma_{xy}$$



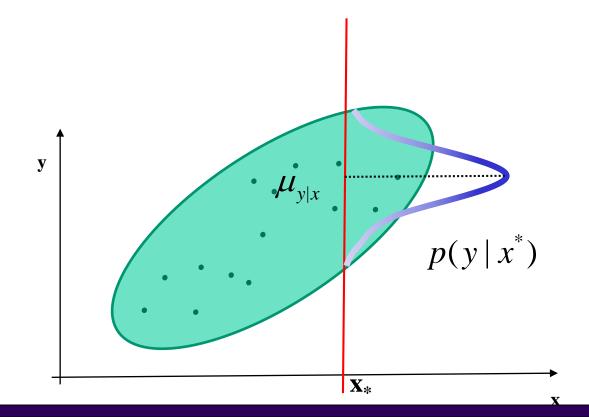




GMR: with a single Gaussian component (k=1)

Given a query point x^* , we compute the conditional.

The mean of the conditional distribution is the model's prediction. $\hat{y}(x^*) = E\{p(y|x^*)\} = \mu_{y|x^*}$



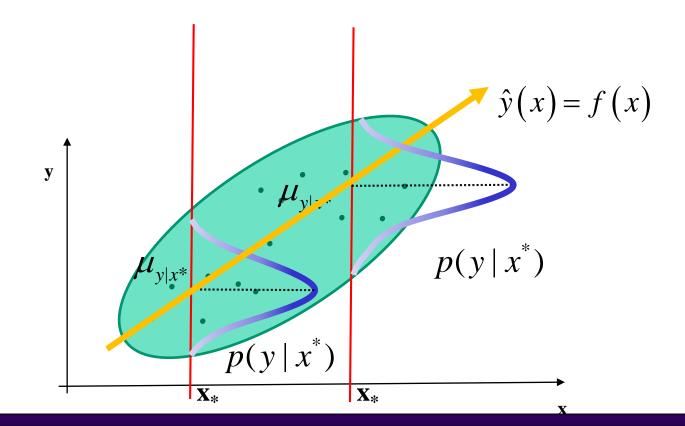




GMR: with a single Gaussian component (k=1)

Given a query point x^* , we compute the conditional.

The mean of the conditional distribution is the model's prediction. $\hat{y}(x^*) = E\{p(y \mid x^*)\} = \mu_{y\mid x^*}$





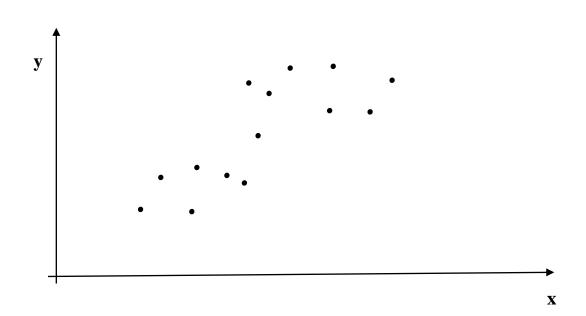


GMR: with multiple Gaussian components (k>1)

1) Estimate the *joint* density, p(x,y), across pairs of datapoints using GMM.

$$p(x,y) = \sum_{k=1}^{K} \alpha_k \cdot p(x,y;\mu^k, \Sigma^k), \quad \text{with } p(x,y;\mu^k, \Sigma^k) = N(\mu^k, \Sigma^k)$$

 μ^i, Σ^i : mean and covariance matrix of Gaussian k.





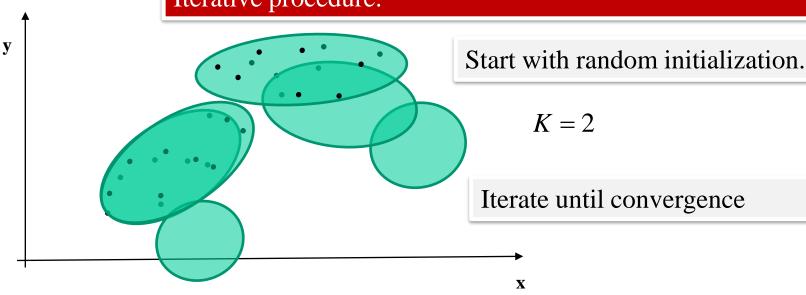
GMR: with multiple Gaussian components (k>1)

1) Estimate the *joint* density, p(x,y), across pairs of datapoints using GMM.

$$p(x,y) = \sum_{k=1}^{K} \alpha_k \cdot p(x,y;\mu^k, \Sigma^k) \qquad \text{with } p(x,y;\mu^k, \Sigma^k) = N(\mu^k, \Sigma^k)$$

 μ^i, Σ^i : mean and covariance matrix of Gaussian k.

Parameters are learned through Expectation-maximization. Iterative procedure.





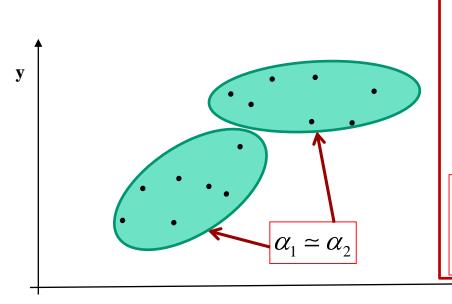


Estimating the joint distribution

1) Estimate the *joint* density, p(x,y), across pairs of datapoints using GMM.

$$p(x,y) = \sum_{k=1}^{K} \alpha_k p(x,y;\mu^k, \Sigma^k), \quad \text{with } p(x,y;\mu^k, \Sigma^k) = N(\mu^k, \Sigma^k)$$

 μ^i, Σ^i : mean and covariance matrix of Gaussian k.



Mixing Coefficients
$$\sum_{k=1}^{R} \alpha_k = 1$$

Relative importance of each Gaussian *k* (*measure how well the Gaussian explains the dataset*):

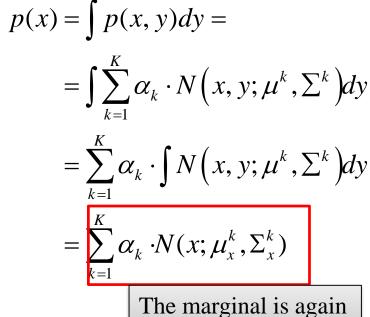
$$\alpha_k = p(k) \cong \frac{1}{M} \sum_{i=1}^{M} \frac{p(x^i, y^i; \mu^k, \Sigma^k)}{\sum_{l} p(x^i, y^i; \mu^l, \Sigma^l)}$$



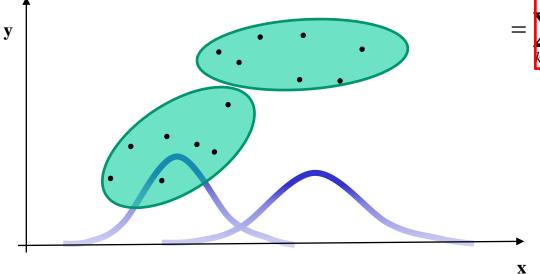
Estimating the marginal distribution

2) Estimate conditional
$$p(y|x) = \frac{p(x,y)}{p(x)}$$

Marginal distribution



a GMM







2) Estimate conditional $p(y|x) = \frac{p(x,y)}{p(x)}$

$$p(x, y) = \sum_{k=1}^{K} \alpha_k \cdot p\left(x, y; \mu^k, \Sigma^k\right) = \sum_{k=1}^{K} \alpha_k \cdot p\left(x; \mu_x^k, \Sigma_x^k\right) p\left(y \mid x; \mu_{y\mid x}^k, \Sigma_{y\mid x}^k\right)$$

The covariance matrix of each
Gauss function k can be
decomposed into blocks of matrices

$$\Sigma^k = \begin{bmatrix} \sum_{xx}^k & \sum_{xy}^k \\ \sum_{yx}^k & \sum_{yy}^k \end{bmatrix},$$

with \sum_{xx}^{k} and \sum_{yy}^{k} the

covariance matrices on x and y and

 \sum_{xy}^{k} the crosscovariance matrix.

$$p(y \mid x) = \frac{\sum_{k=1}^{K} \alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k) p(y \mid x; \mu_{y|x}^k, \Sigma_{y|x}^k)}{\sum_{k=1}^{K} \alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}$$

Set
$$\beta_k(x) = \frac{\alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}{\sum_{k=1}^K \alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}$$





2) Estimate conditional
$$p(y|x) = \frac{p(x,y)}{p(x)}$$

The conditional is again a GMM

$$p(y|x) = \sum_{k=1}^{K} \beta_{k}(x) p(y|x; \mu_{y|x}^{k}, \Sigma_{y|x}^{k})$$

Weight of the marginals

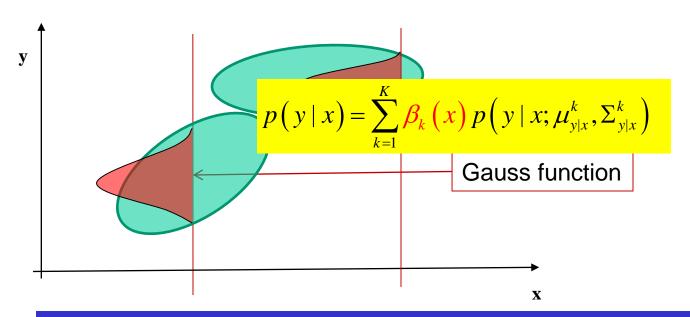
Set
$$\beta_k(x) = \frac{\alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}{\sum_{k=1}^K \alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}$$







$$p(y | x; \boldsymbol{\mu}_{y|x}^k, \boldsymbol{\Sigma}_{y|x}^k)$$

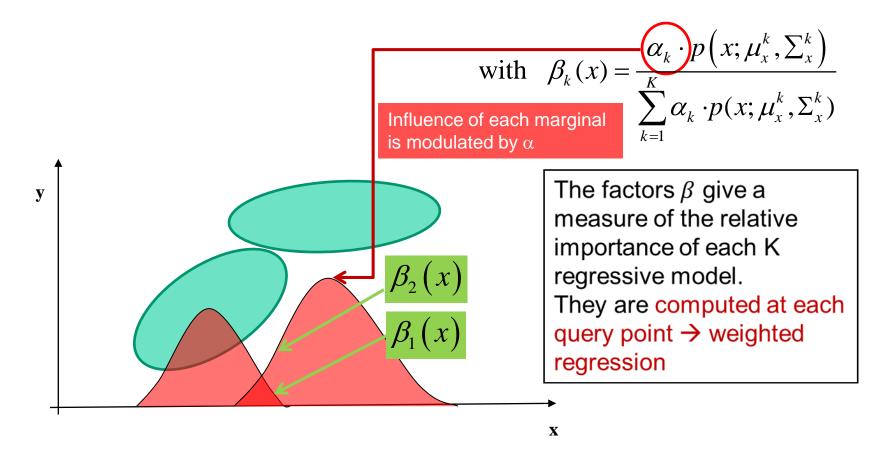


The expression changes depending on the query point





$$p(y \mid x) = \sum_{k=1}^{K} \beta_k(x) p(y \mid x; \mu_{y|x}^k, \Sigma_{y|x}^k)$$





Computing the regression

The GMR prediction at a query point x^* is obtained by computing $E\{p(y/x)\}$

$$y = E\{p(y|x)\} = \sum_{k=1}^{K} \beta_{k}(x) \left(\mu_{y}^{k} + \sum_{yx}^{k} \left(\sum_{xx}^{k}\right)^{-1} \left(x - \mu_{x}^{k}\right)\right)$$

$$\Rightarrow y = \sum_{k=1}^{K} \beta_{k}(x) \cdot \tilde{\mu}_{y|x}^{k}(x)$$
Linear combination of K local regressive models
$$\tilde{\mu}_{y|x}^{2}(x)^{y}$$

$$\tilde{\mu}_{y|x}^{1}(x)$$

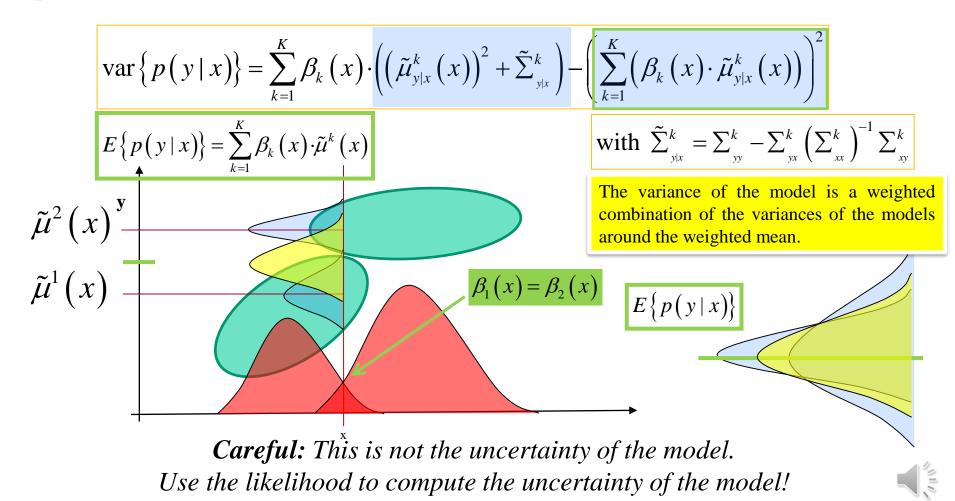
$$\beta_{2}(x)$$

$$\beta_{1}(x)$$



Computing the variance

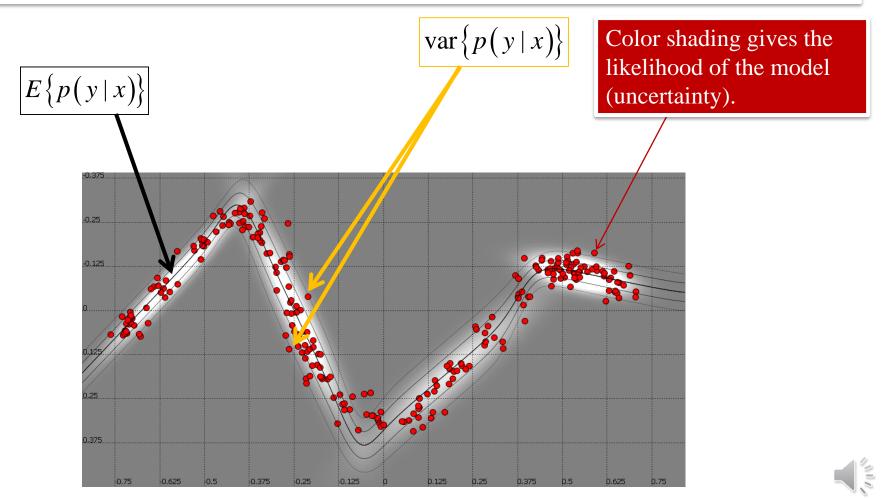
Computing the variance $var\{p(y|x)\}$ provides information on the <u>variability of the prediction</u>.



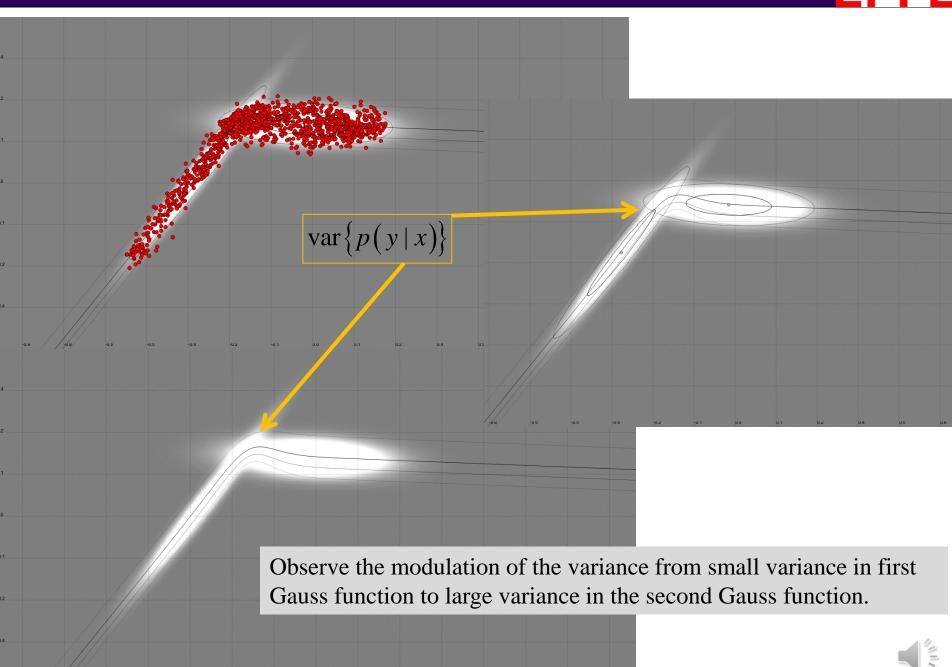


Interpreting the variance

The variance $var\{p(x,y)\}$ can be visualized all around the regressive line. It provides information on the evolution of the noise in state space.



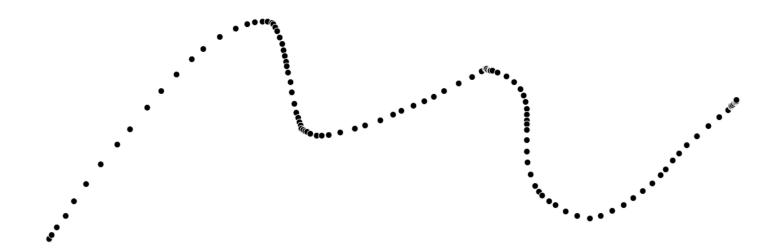




19

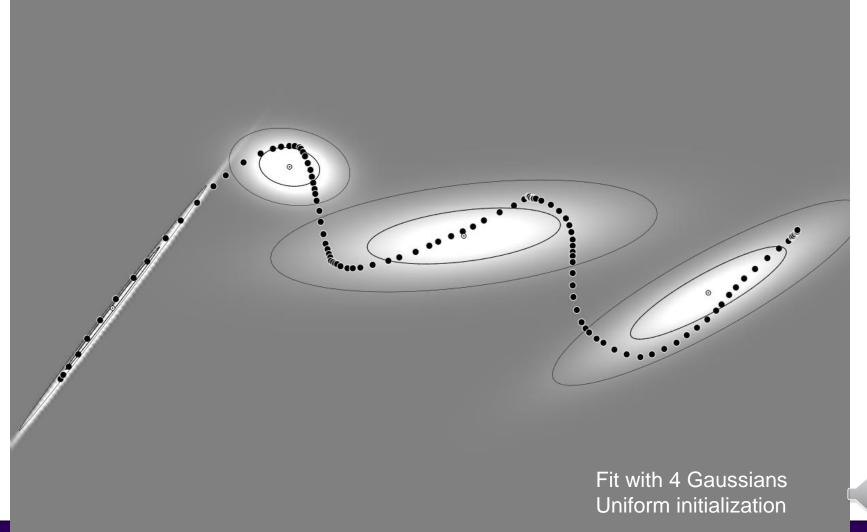


GMR: Sensitivity to Choice of *K* and Initialization



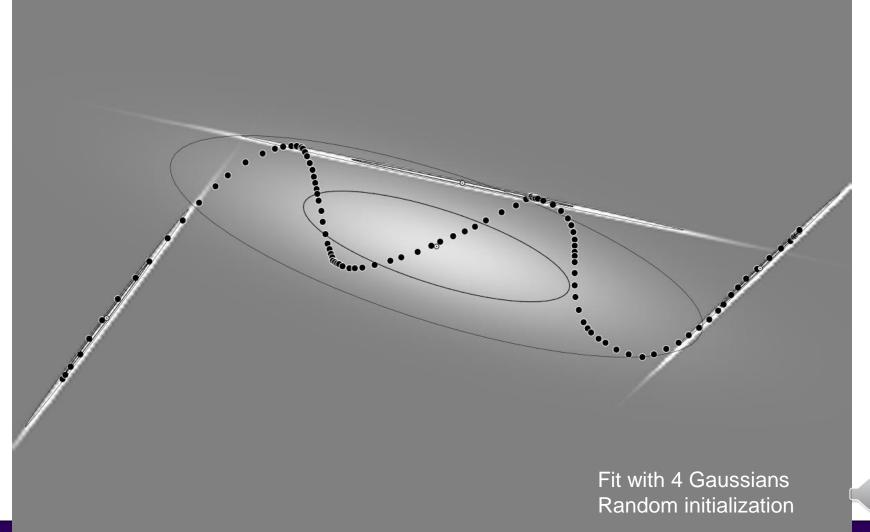


GMR: Sensitivity to Choice of *K* and Initialization





GMR: Sensitivity to Choice of *K* and Initialization





GMR: Take home message

GMR proceeds in two steps:

- \square Parametrize the density p(x,y) and then estimate solely the parameters. The density is constructed from a mixture of K Gaussians
- ☐ Compute the regression from the expectation over the conditional.

Such *generative model* provides more information than models that directly computing p(y|x).

- → It allows to learn to predict a multi-dimensional output y.
- → It allows to embed correlations across the output dimensions y.
- \rightarrow It allows to query x given y, i.e. to compute p(x/y).

It comes at the cost of computing the full distribution, while often we care only for the conditional.

The estimate depends on E-M which is very sensitive to initialization.