

December 19

Lecture on site and online (zoom)

Video of lecture will be made available after the lecture

9h15-12am: Handling Imbalanced Datasets

Enable Incremental Learning

Overview of the course

Exam preparation

•12am-13: Open Q&A regarding exam and material of the class



Q&A session – exam preparation

December 19, 11h15-13:00

Friday 24, 11h00-12h00 room ME.A3.31



Nonlinear Regression

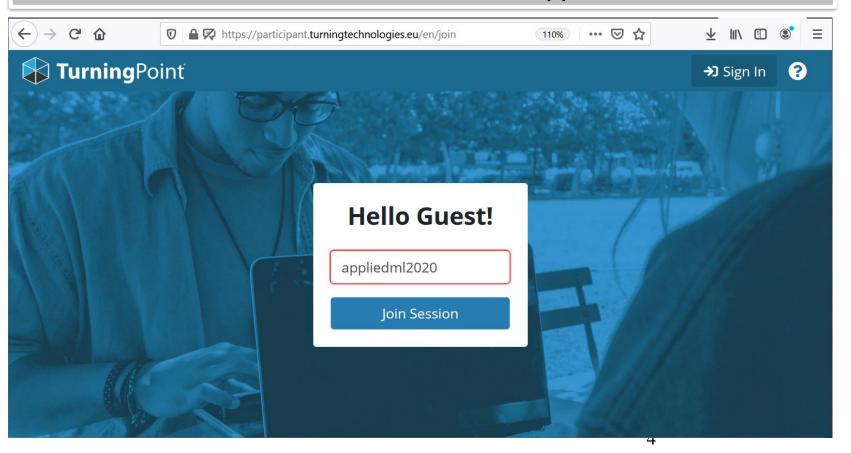
Interactive Lecture



Launch polling system

https://participant.turningtechnologies.eu/en/join

Acces as GUEST and enter the session id: appliedml2020





Linear Regression

Find the optimal parameter w through least-square regression:

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} (w^T x^i - y^i)^2 \right)$$

Closed-form solution:

$$w^* = \left(XX^T\right)^{-1} Xy$$

Weighted Regression

Find the optimal parameter w through

least-square regression:

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2 \right)$$

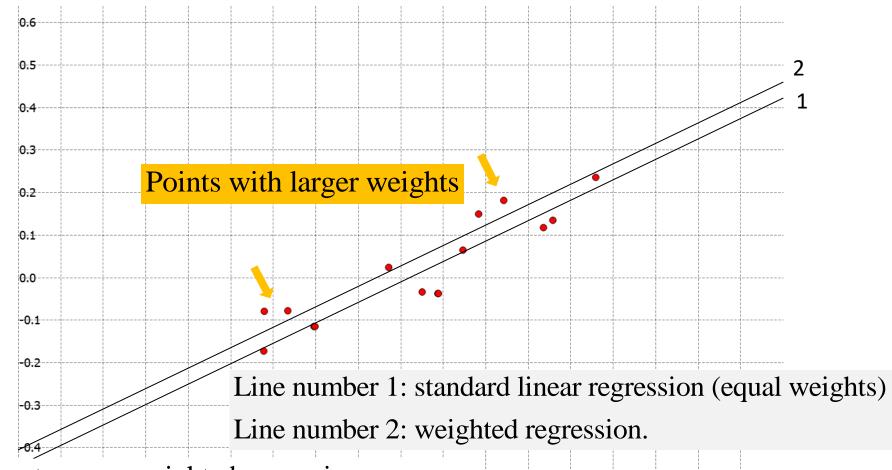
Closed-form solution:

$$w^* = \left(ZZ^T\right)^{-1} Zv x$$

$$Z = XB^{1/2}$$
 and $v = B^{1/2}y$



Weighted linear regression



Least-square weighted regression

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2 \right), \quad \beta_i > 0 \quad \beta_1 = \beta_2 \dots = \beta_M$$

$$\beta_i > 0$$
 $\beta_1 = \beta_2 ... = \beta_M$ 0.6 0.7 0.8



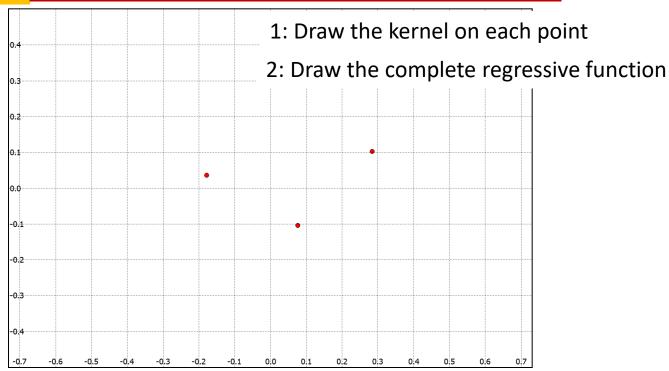
Locally weighted regression

Introduce local solution:

$$\widehat{y}(x) = \sum_{i=1}^{M} \beta_i(x) y^i / \sum_{j=1}^{M} \beta_j(x) \qquad \beta_i(x) \in \mathbb{R} : \text{ weights function of } x$$

$$\beta_i(x) = e^{-\|x^i - x\|^2} \text{ Also closed-form solution, but local regression}$$

$$\beta_i(x) = e^{-\|x^i - x\|^2}$$
 Also closed-form solution, but local regression





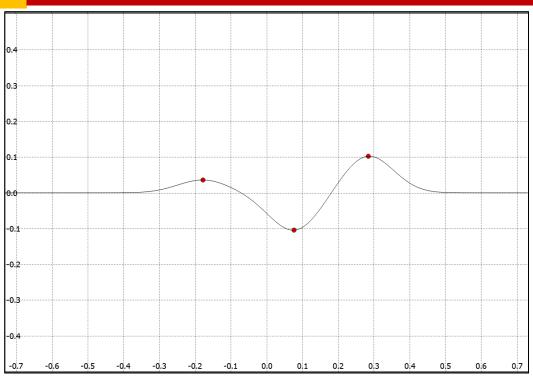
Locally weighted regression

Introduce local solution:

$$\widehat{y}(x) = \sum_{i=1}^{M} \beta_i(x) y^i / \sum_{j=1}^{M} \beta_j(x) \quad \beta_i(x) \in \mathbb{R} : \text{ weights function of } x$$

$$\beta_i(x) = e^{-\|x^i - x\|^2} \text{ Also closed-form solution, but local regression}$$

$$\beta_i(x) = e^{-\|x^i - x\|^2}$$
 Also closed-form solution, but local regression

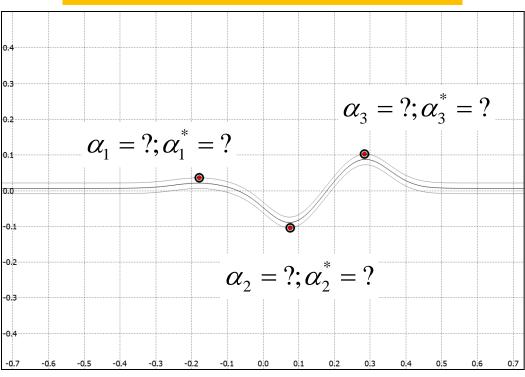




Support Vector Regression

SVR determines automatically which point matters for building the regression.

$$y = f(x) = \sum_{i=1}^{M} (\alpha_i - \alpha_i^*) k(x^i, x) + b$$

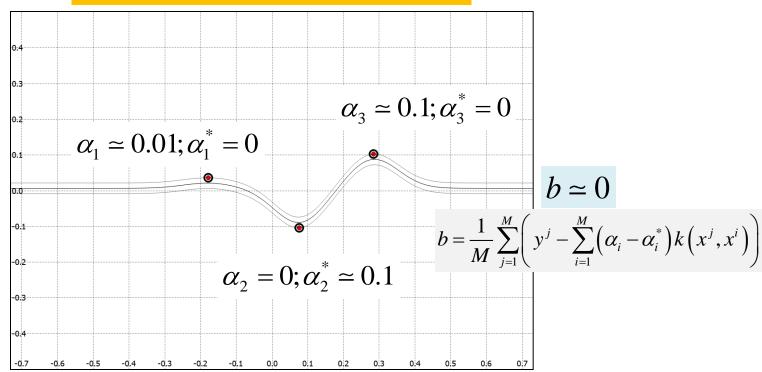




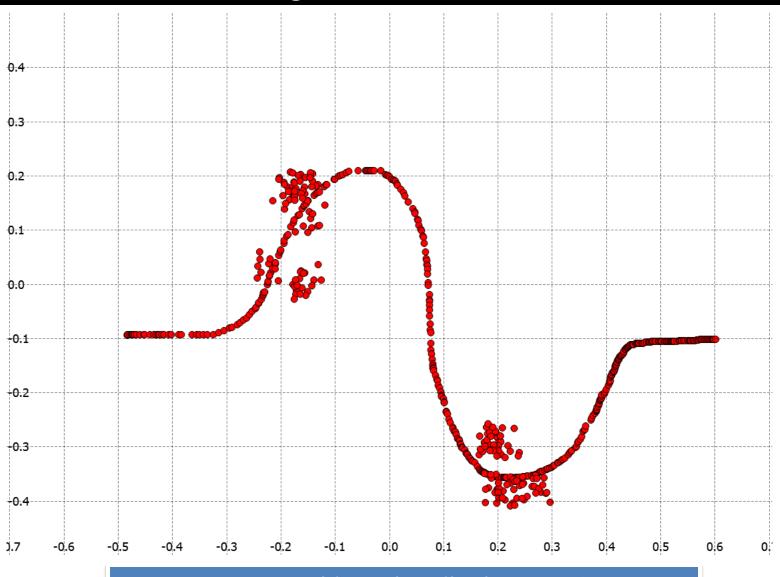
Support Vector Regression

Determines automatically which point matters for building the regression.

$$y = f(x) = \sum_{i=1}^{M} (\alpha_i - \alpha_i^*) k(x^i, x) + b$$

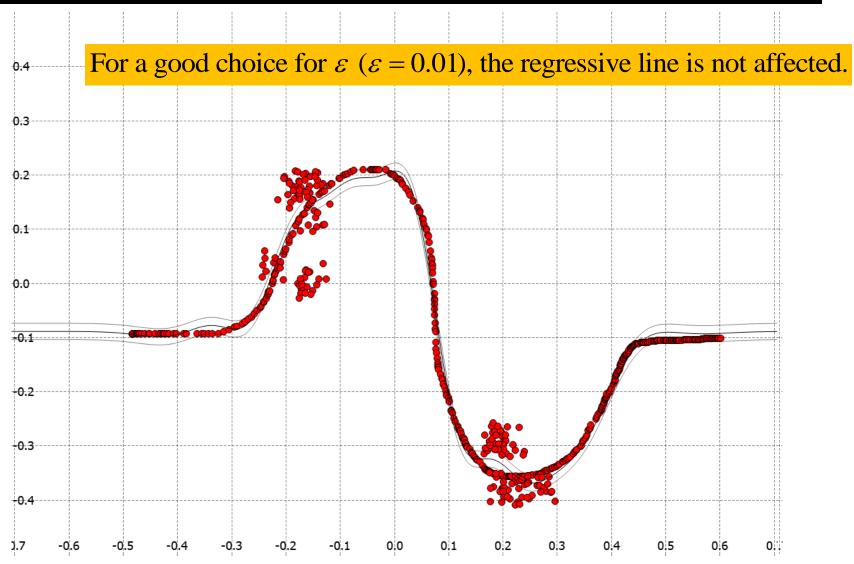




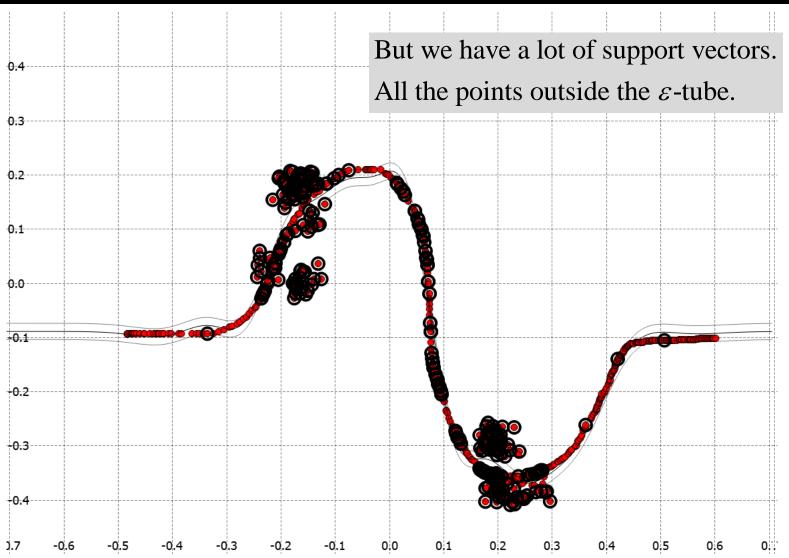


How would SVR handle this noise?

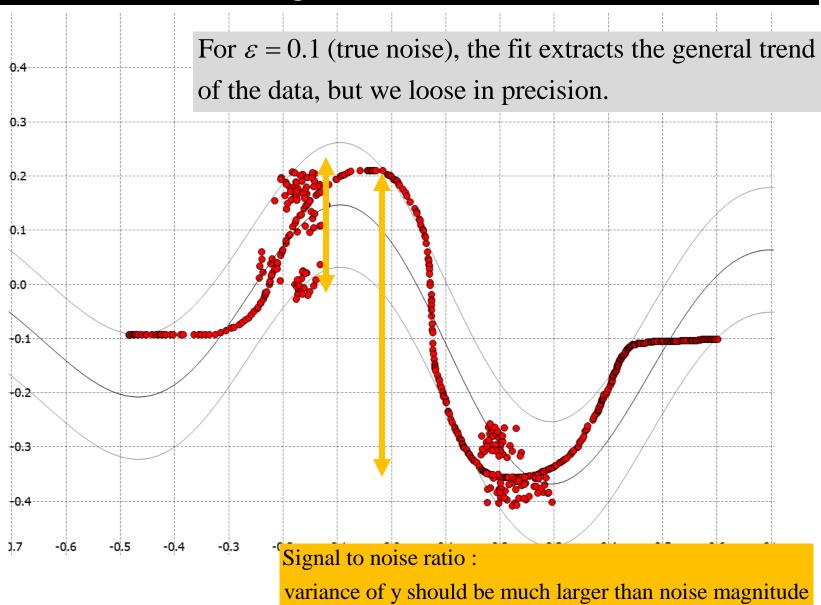




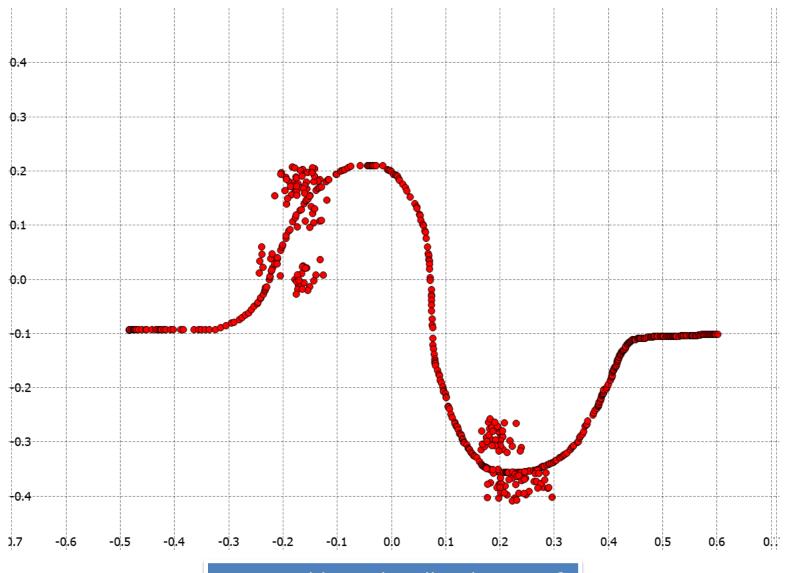






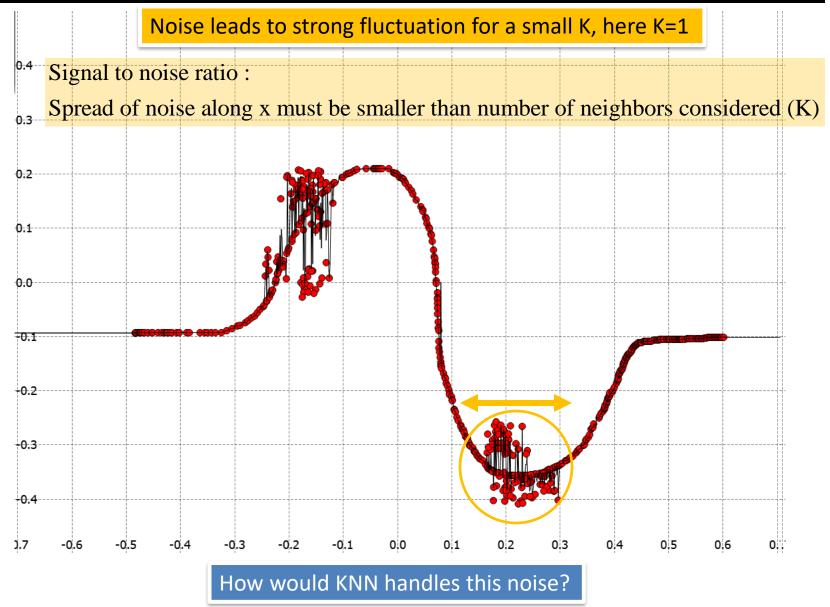






How would KNN handles this noise?

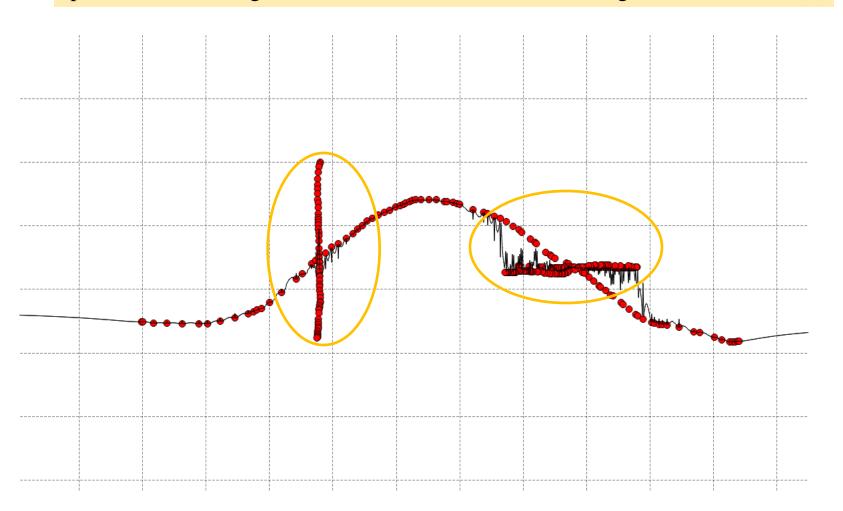




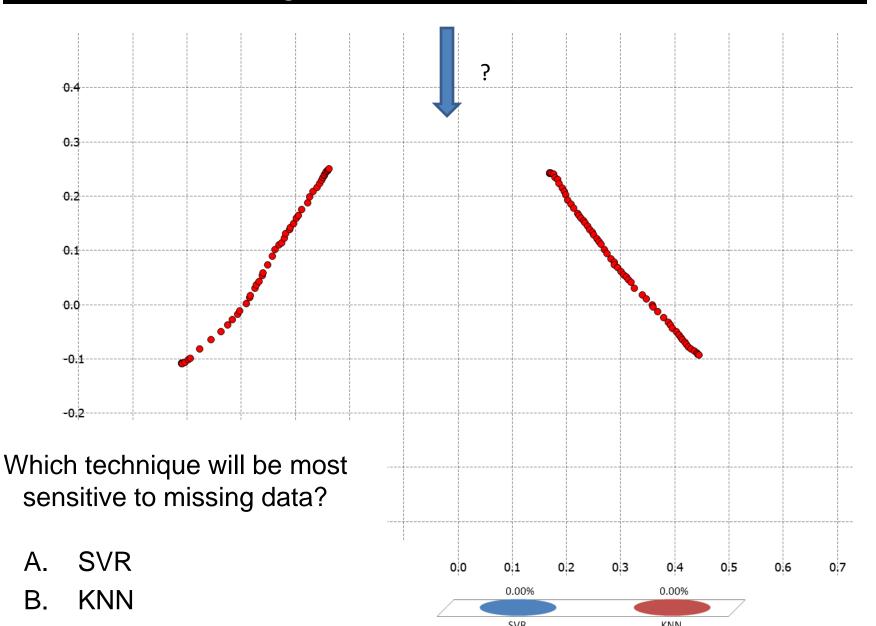


Signal to noise ratio:

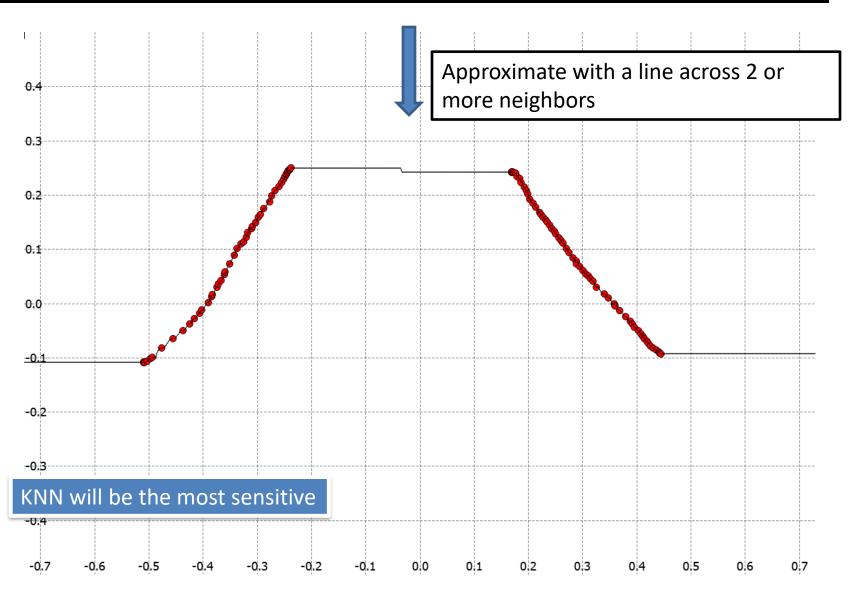
Spread of noise along x must be smaller than number of neighbors considered (K)



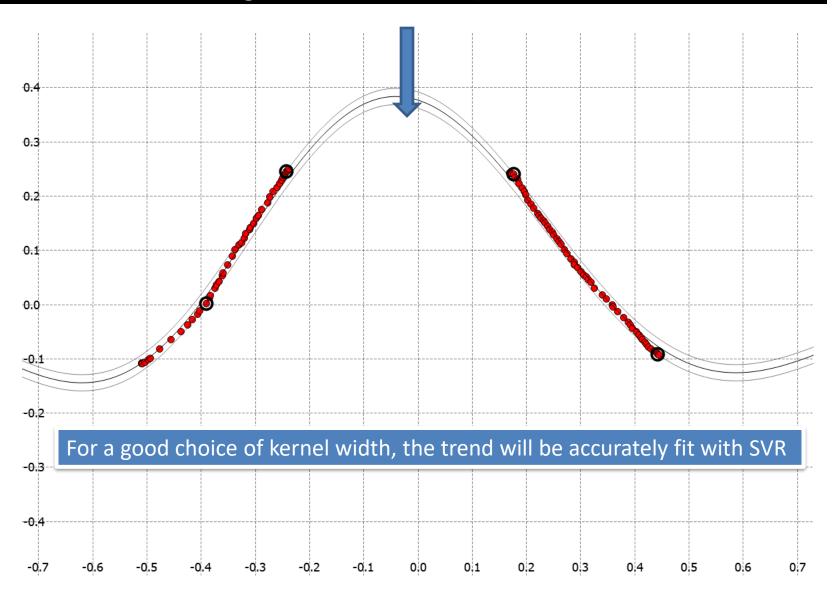




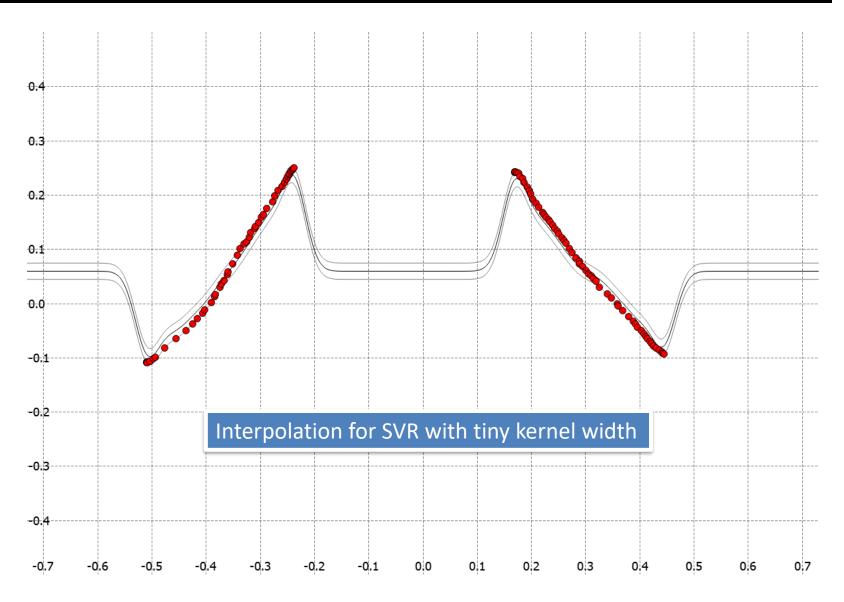








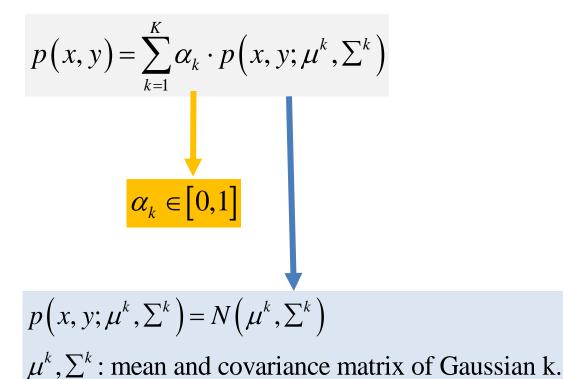






Gaussian mixture regression

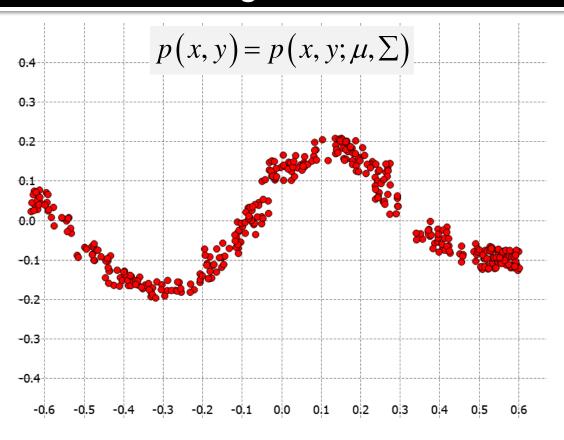
We must first learn the joint distribution



And then we compute the regressive signal:

$$y = E\{p(y \mid x)\}$$







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

$$K = 1 \Rightarrow y = \mu_{y} + \sum_{yx} \left(\sum_{xx}\right)^{-1} \left(x - \mu_{x}\right)$$

$$y = \sum_{yx} \left(\sum_{xx}\right)^{-1} x + \mu_{y} - \sum_{yx} \left(\sum_{xx}\right)^{-1} \mu_{x}$$
with $\beta_{k}(x) = \frac{\alpha_{k} \cdot p\left(x; \mu_{x}^{k}, \sum_{x}^{k}\right)}{\sum_{k=1}^{K} \alpha_{k} \cdot p\left(x; \mu_{x}^{k}, \sum_{x}^{k}\right)}$

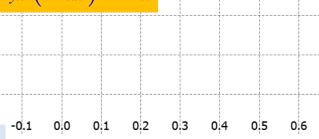
For which model is the regressive curve a straight line?

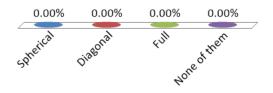


B. Diagonal

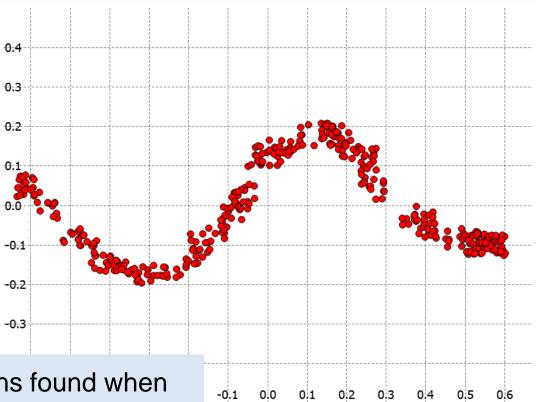
C. Full

D. None of them







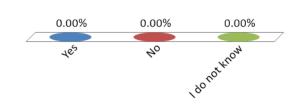


Are the solutions found when using spherical or diagonal covariance matrices different?



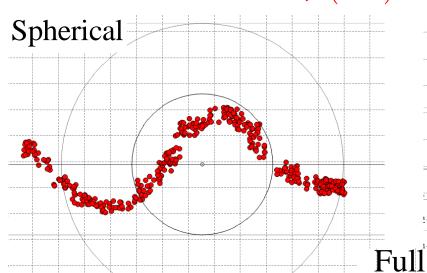
B. No 🤳

C. I do not know









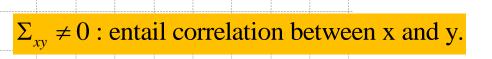
Diagonal $\Sigma_x = \begin{bmatrix} \sigma_{x_1} & 0 \\ 0 & \sigma_{x_2} \end{bmatrix}, \Sigma_y = \begin{bmatrix} \sigma_{y_1} & 0 \\ 0 & \sigma_{y_2} \end{bmatrix}, \Sigma_{xy} = 0$

x, y unidimensional, hence Σ_x , Σ_y are identical to spherical case.



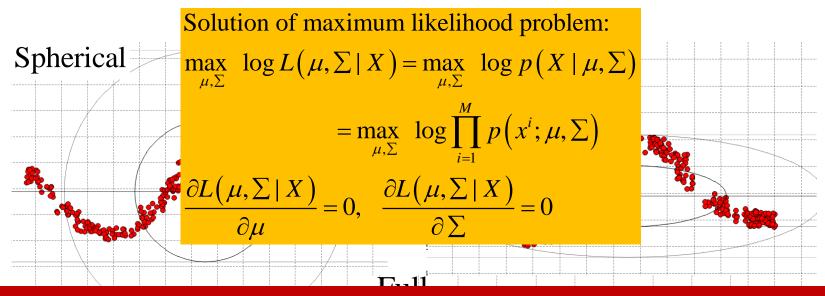
 Σ_{xy} is a matrix for multidimensional x, y.

 Σ_{xy} entails the correlation across x and y.





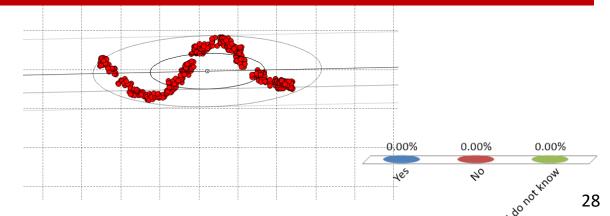
GMR – 1 Gauss fct: Uniqueness of solution



The solution is unique and found in closed-form: mean and variance of dataset, see exercises of pdf lecture

Is the solution unique?

- A. Yes
- B. No
- C. I do not know

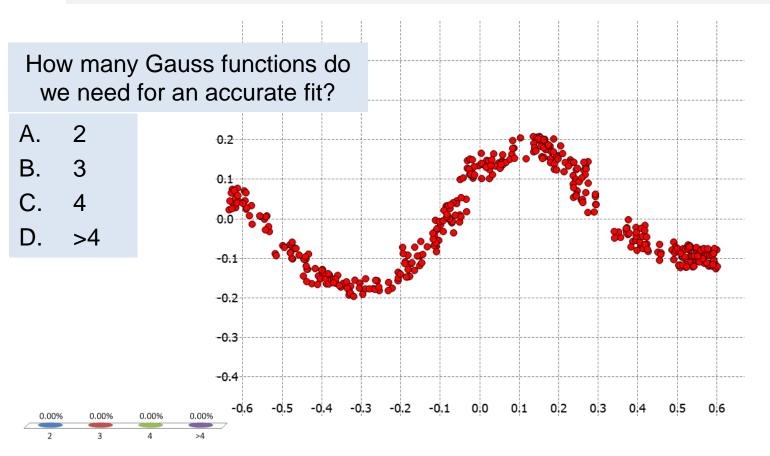




GMR: Multiple Gauss Functions

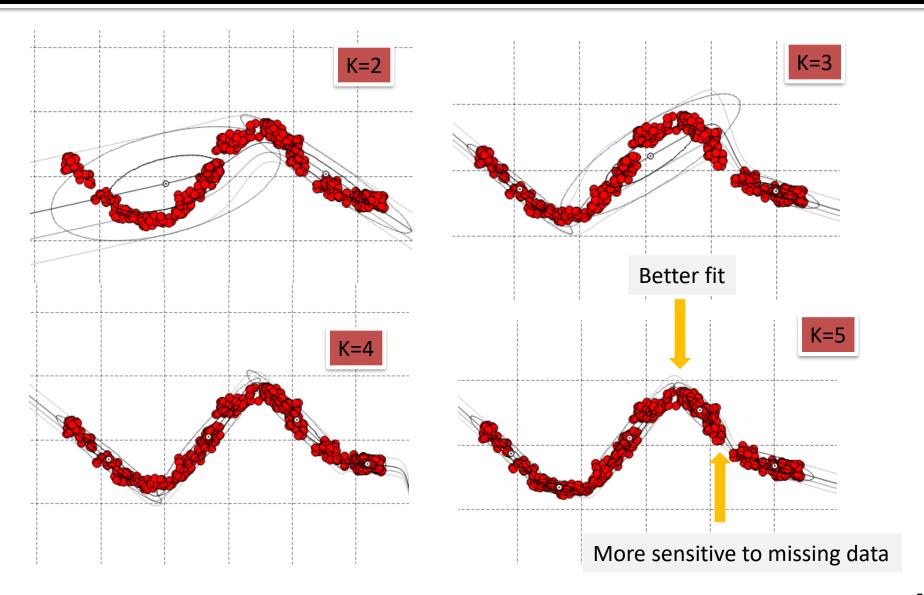
$$p(x,y) = \sum_{k=1}^{K} \alpha_k \cdot p(x,y;\mu^k, \Sigma^k), \quad \text{with } p(x,y;\mu^k, \Sigma^k) = N(\mu^k, \Sigma^k)$$

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

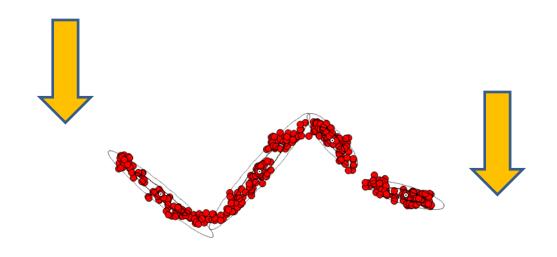




Accuracy of fit with multiple Gauss functions

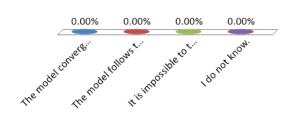


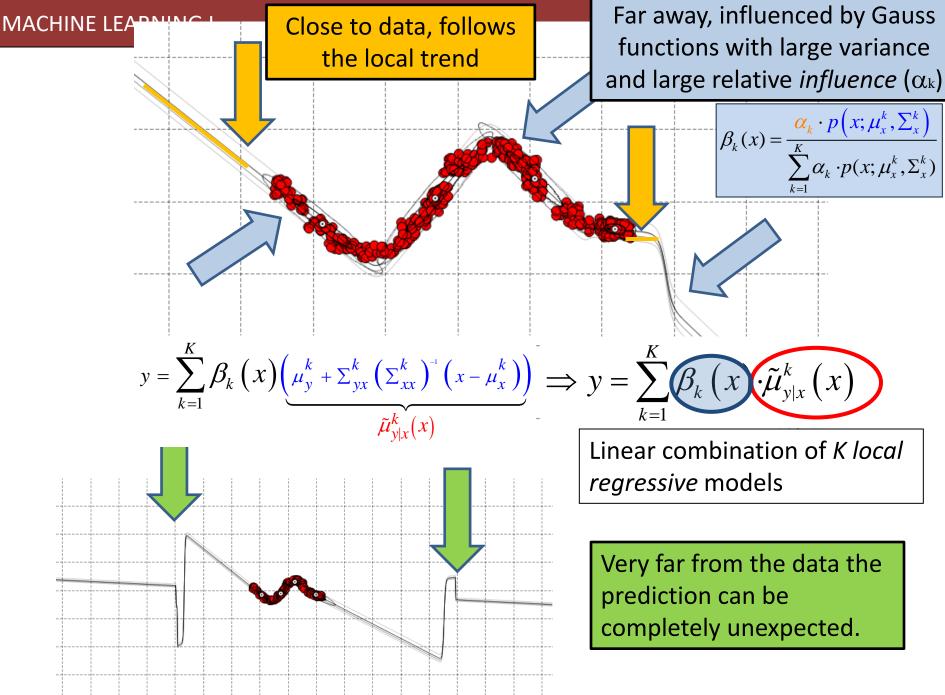




What does the model predict away from the data?

- A. The model converges to a single value.
- B. The model follows the local trend.
- C. It is impossible to tell.
- D. I do not know.



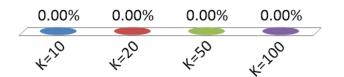


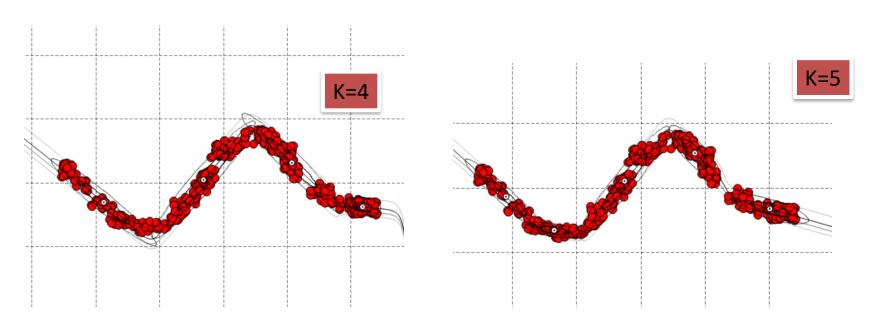


Knowing that we have ~100 points and use a 2/3rd training/testing ratio, for which K would we start seeing overfitting?



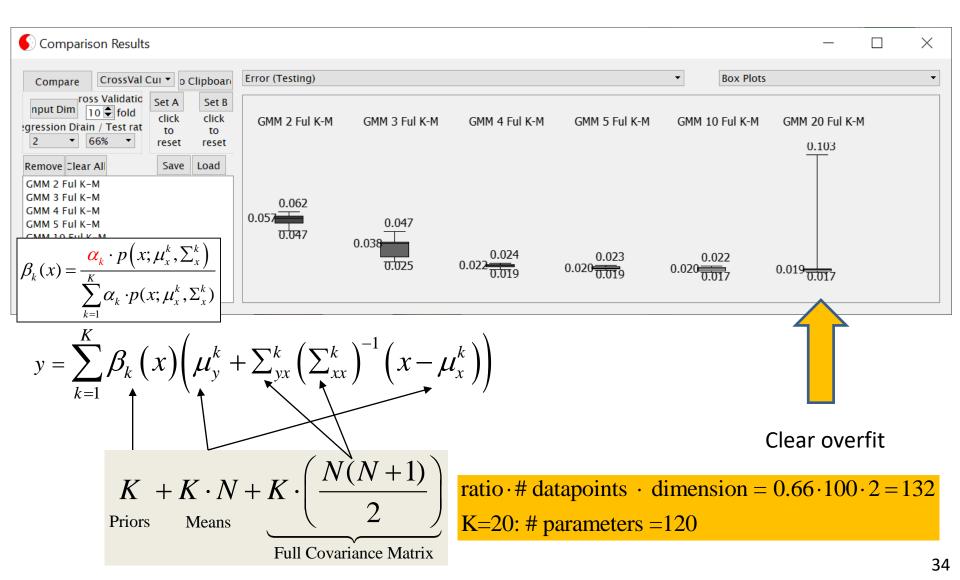
D. K=100







Overfitting with multiple Gauss functions

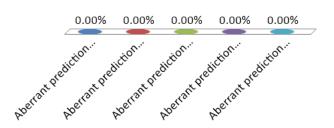




What would be the effect of overfitting in GMR?

Multiple correct responses

- A. Aberrant prediction when far from the dataset
- B. Aberrant prediction even for query points close to the dataset
- C. Aberrant prediction could be any value, even values never seen at training.
- D. Aberrant prediction would be a value that remains within variance of the dataset.
- E. Aberrant prediction can only be "zero".



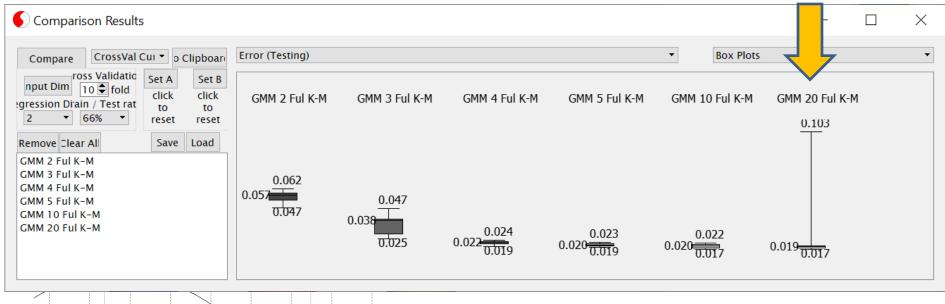


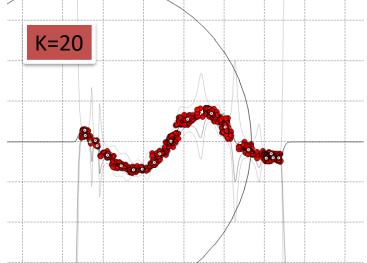
Overfitting with multiple Gauss functions

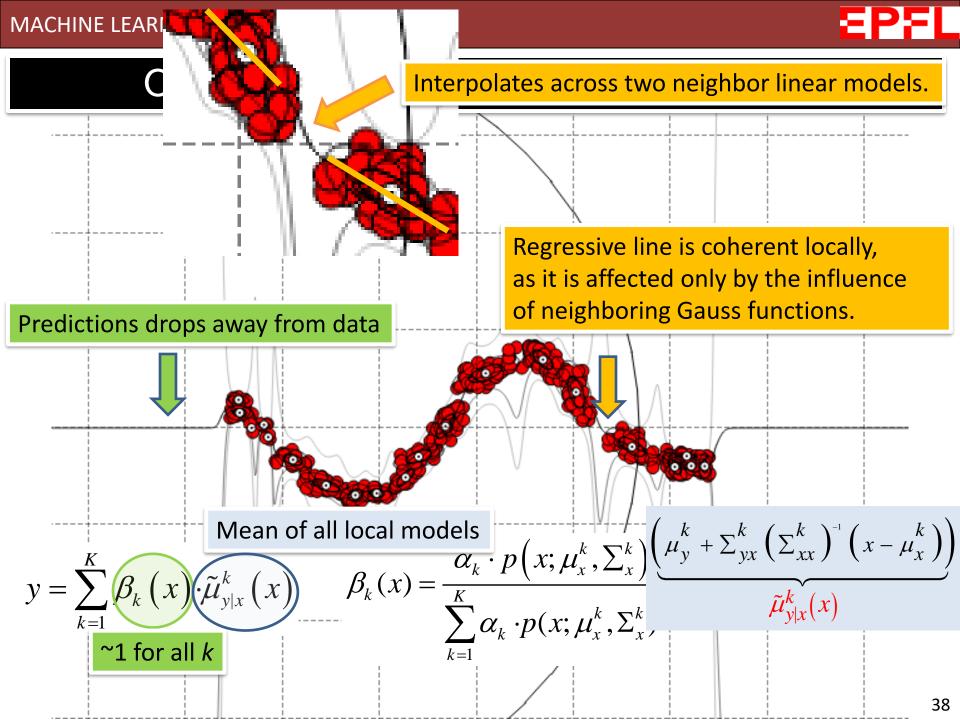




Overfitting with multiple Gauss functions



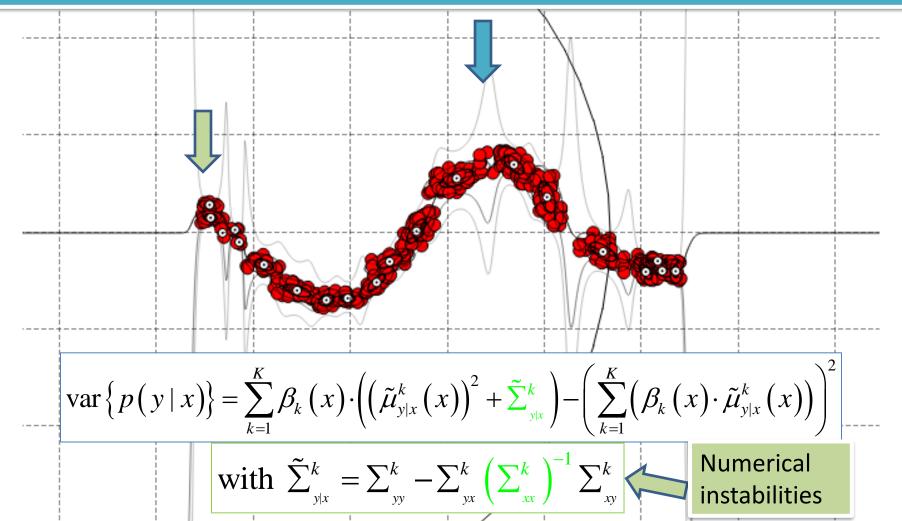






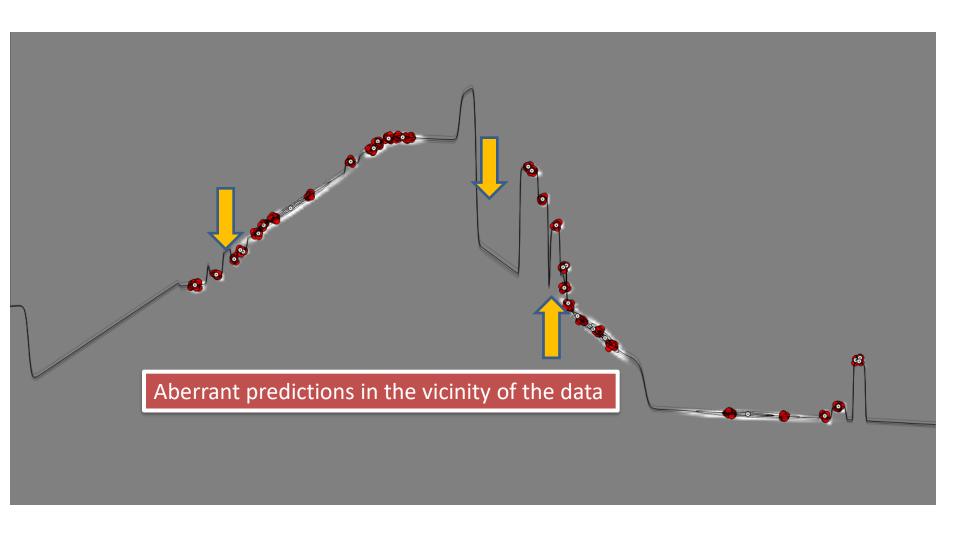
Overfitting with multiple Gauss functions

Aberrant predictions of variance –not enough statistics to estimate all parameters, when single full Gauss function estimated from too few datapoints.



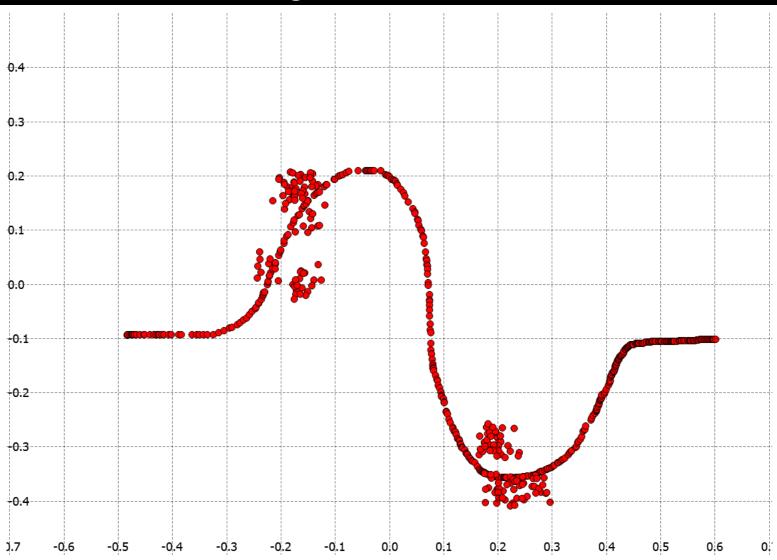


Overfitting with multiple Gauss functions





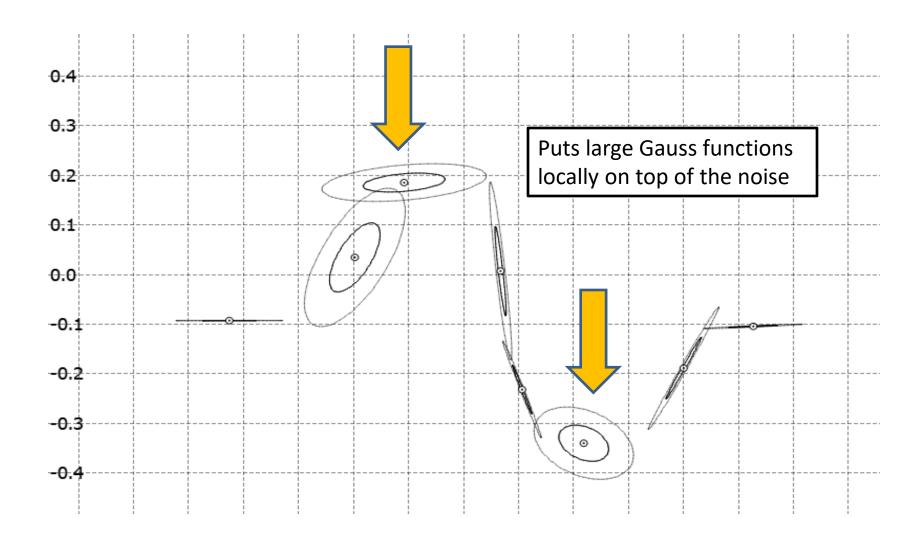
Regression: noise



How would GMR handle this noise?

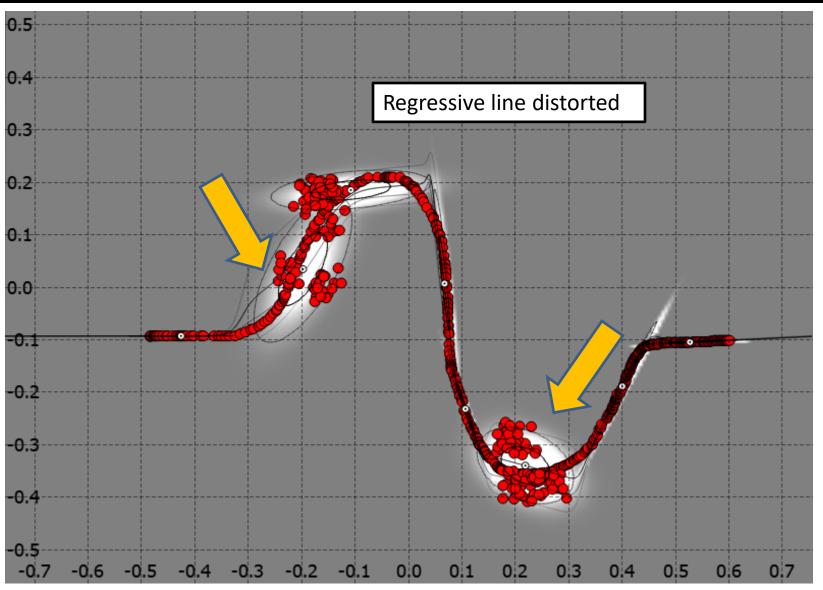


Regression: noise



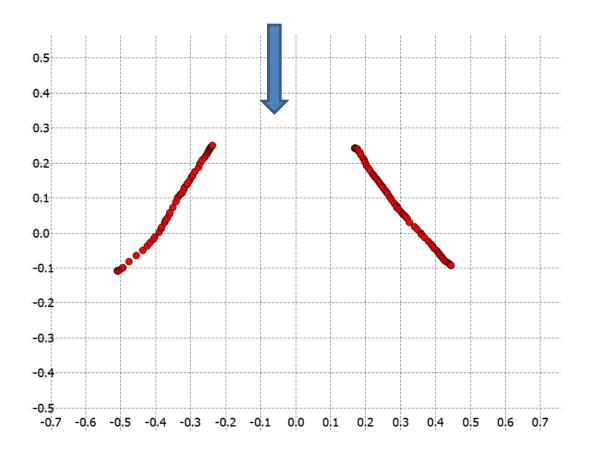


Regression: noise





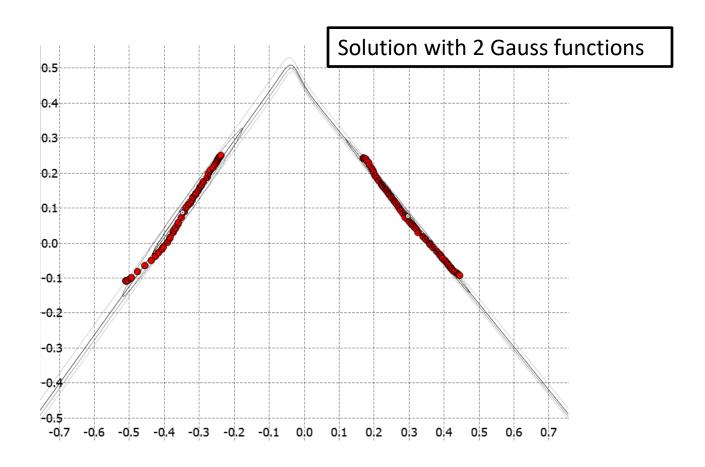
Regression: interpolation



How would GMR handle missing data?

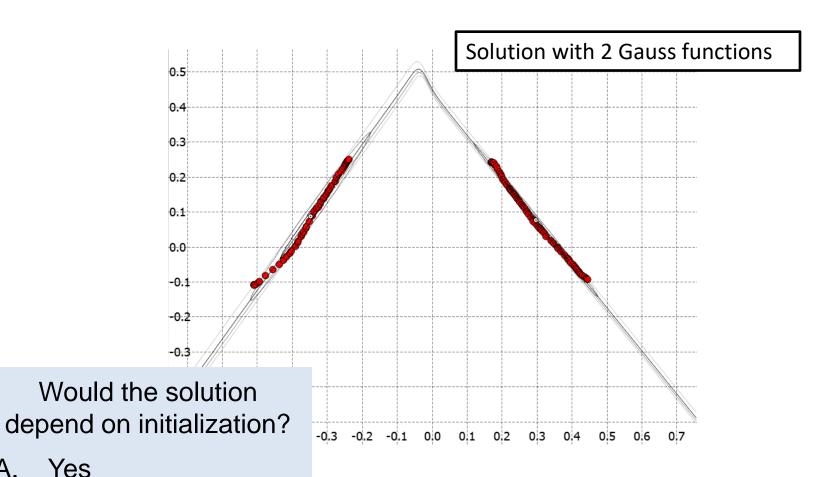


Regression: interpolation



GMR interpolates correctly following the trend with a small curvature at the junction

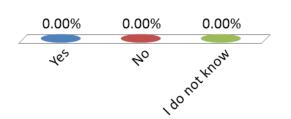




Yes Α.

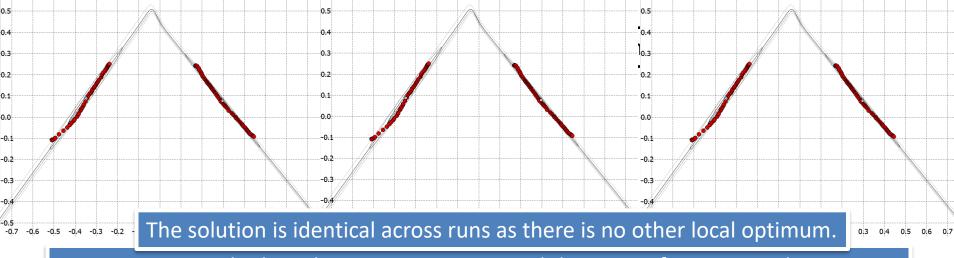
B. No

C. I do not know

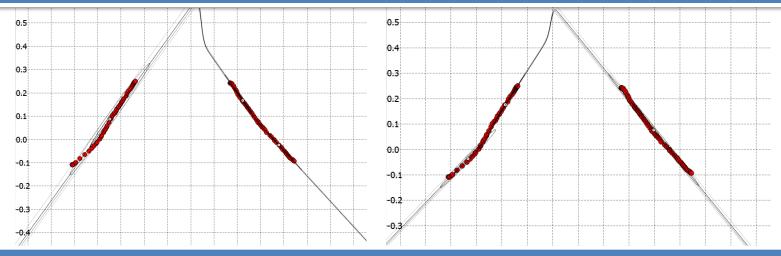




Regression: interpolation

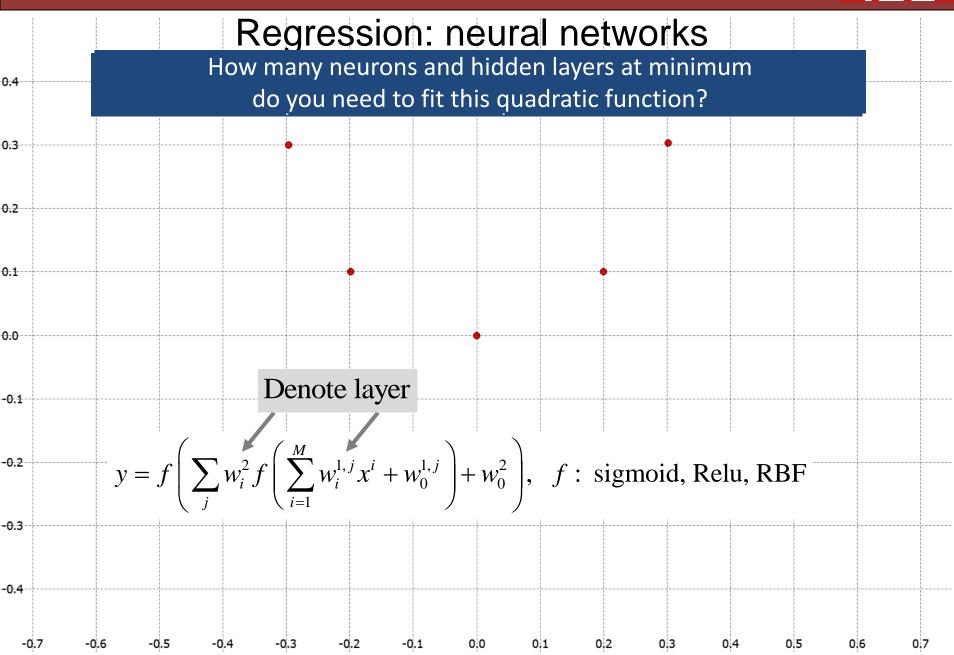


But in principle the solution is not unique and there are often many solutions, each of which corresponds to a local optima on the likelihood.



For instance, we find 2 distinct solutions for a GMM with K=3.

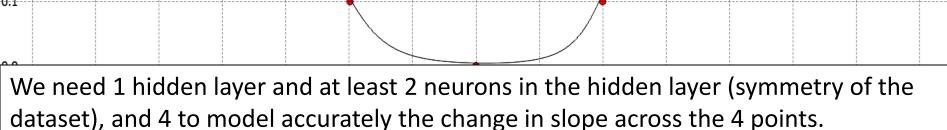


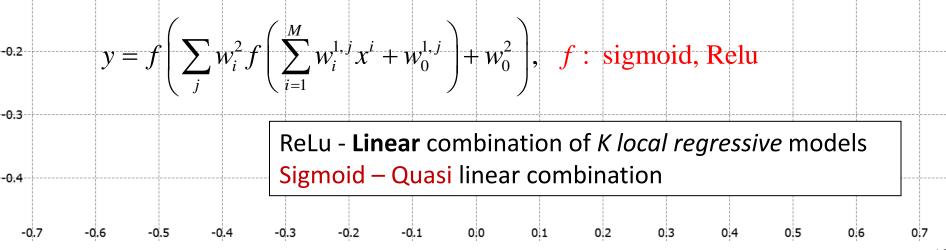




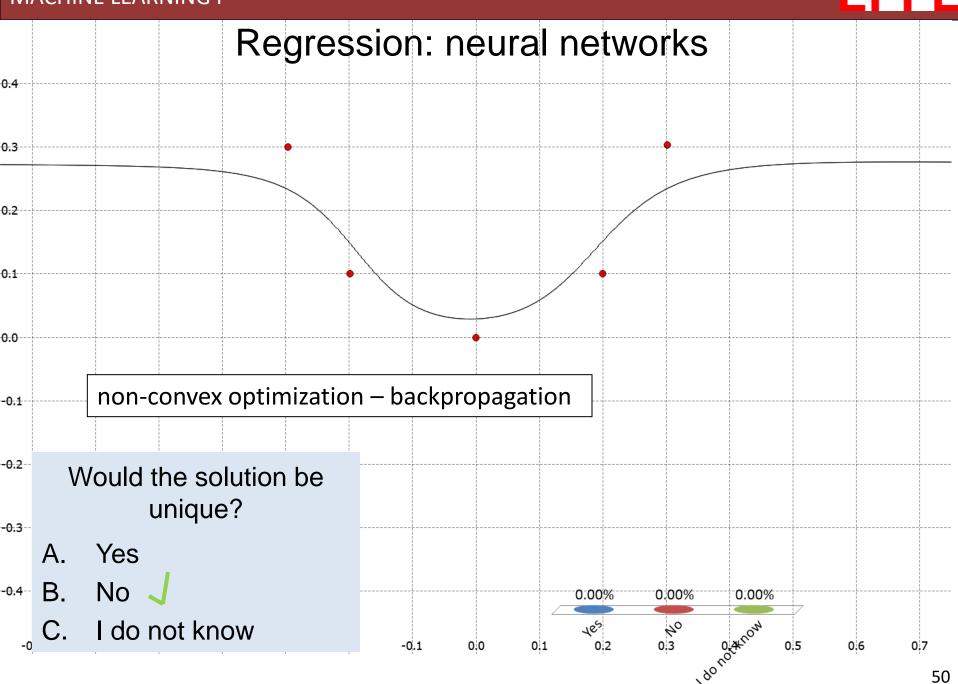


How many neurons and hidden layers at minimum do you need to fit this quadratic function?

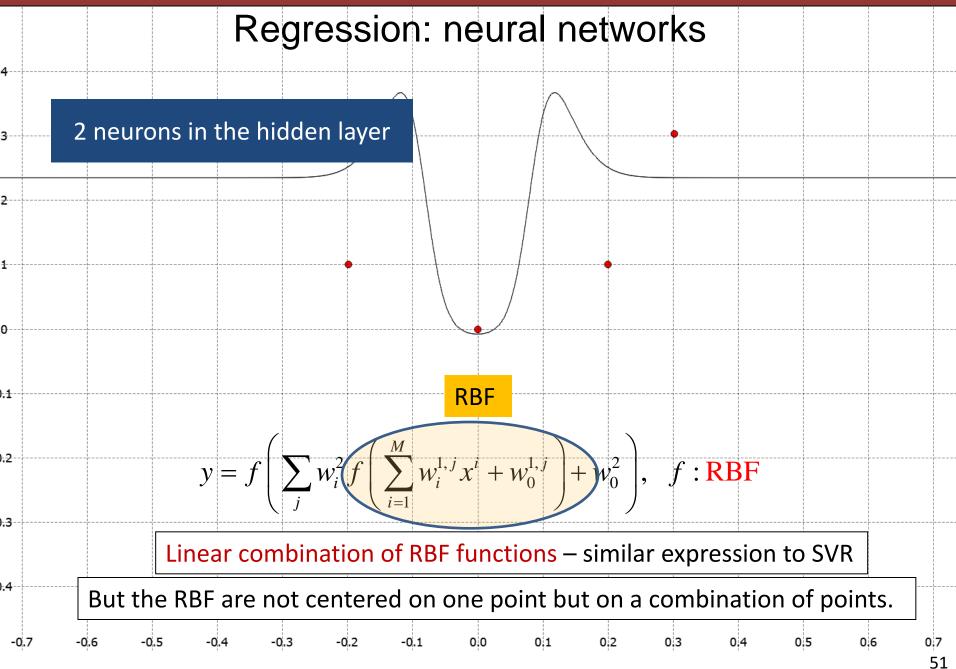




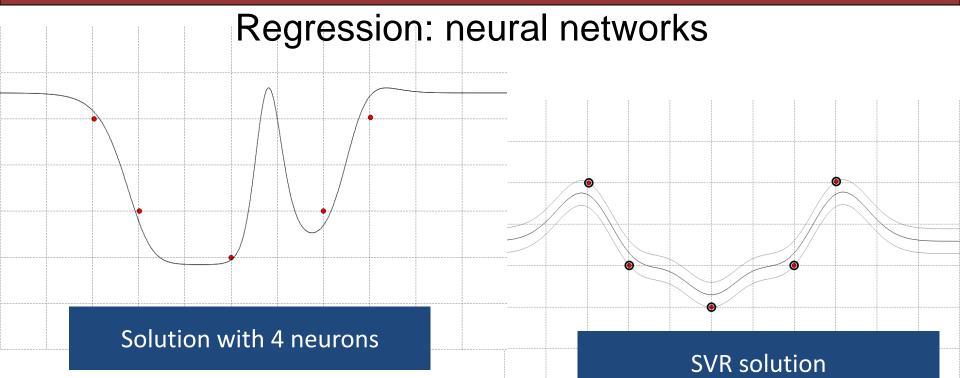












$$y = f\left(\sum_{j} w_{i}^{2} f\left(\sum_{i=1}^{M} w_{i}^{1,j} x^{i} + w_{0}^{1,j}\right) + w_{0}^{2}\right), \quad f : \mathbf{RBF}$$

RBF function – similar expression to SVR

But non-convex optimization – backpropagation, in contrast to SVR