

December 18

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

January 26, 12h00-14:00

Room ME.A3.31

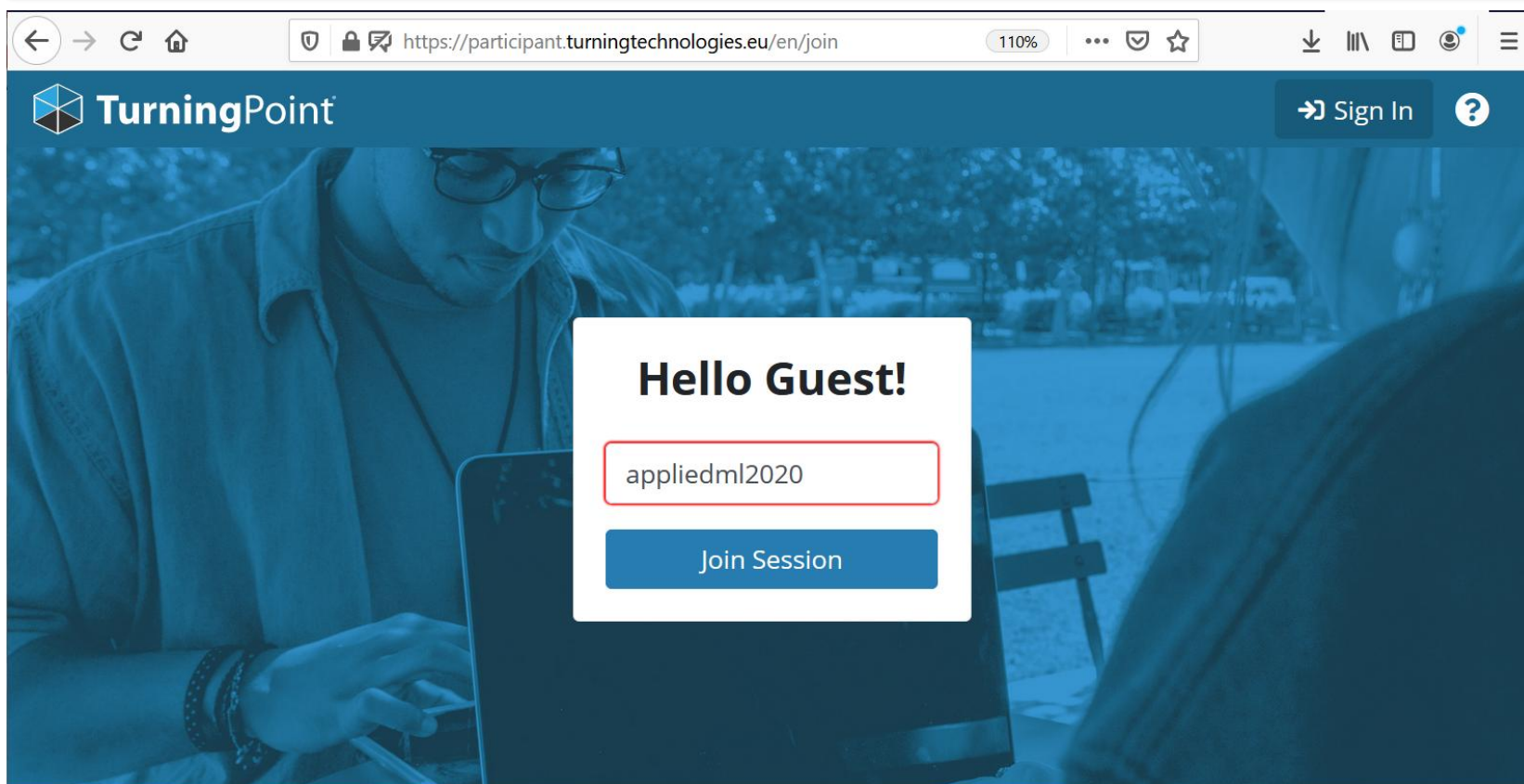
Nonlinear Regression

Interactive Lecture

Launch polling system

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

Access as GUEST and enter the session id: *appliedml2020*



The screenshot shows a web browser window displaying the TurningPoint participant interface. The browser's address bar shows the URL <https://participant.turningtechnologies.eu/en/join>. The page features the TurningPoint logo in the top left and a 'Sign In' button in the top right. A central dialog box titled 'Hello Guest!' contains a text input field with the session ID 'appliedml2020' and a blue 'Join Session' button below it. The background of the page is a blue-tinted image of a person using a laptop.

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:

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

Weighted Regression

Find the optimal parameter w through
least-square regression:

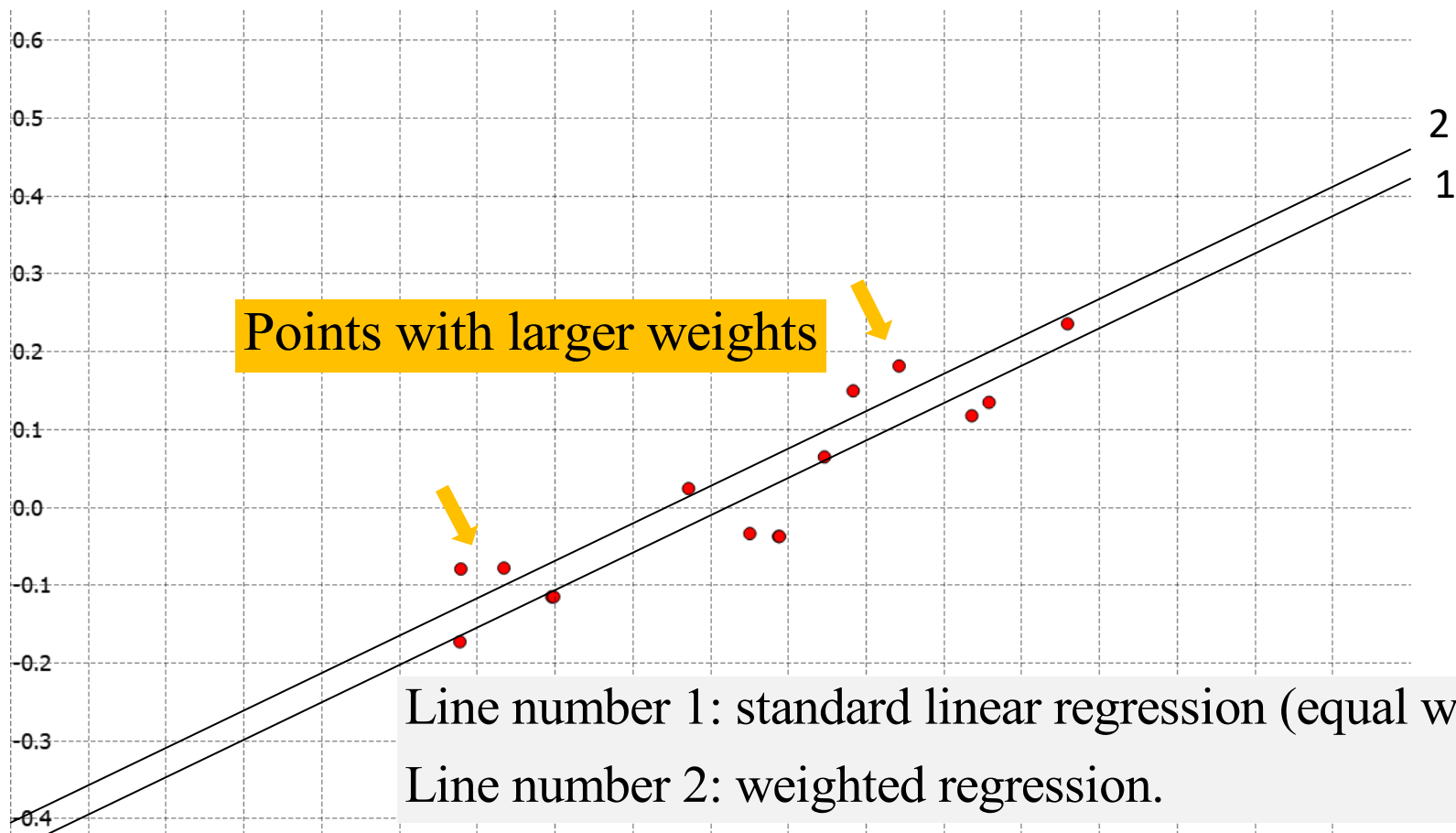
$$w^* = \min_w \left(\sum_{i=1}^M \frac{1}{2} \beta_i (w^T x^i - y^i)^2 \right)$$

Closed-form solution:

$$w^* = (ZZ^T)^{-1} Zv$$

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

Weighted linear regression



Least-square weighted regression

$$w^* = \min_w \left(\sum_{i=1}^M \frac{1}{2} \beta_i (w^T x^i - y^i)^2 \right), \quad \beta_i > 0 \quad \underbrace{\beta_1 = \beta_2 \dots = \beta_M}_{0.6 \quad 0.7 \quad 0.8}$$

→ Standard linear regression

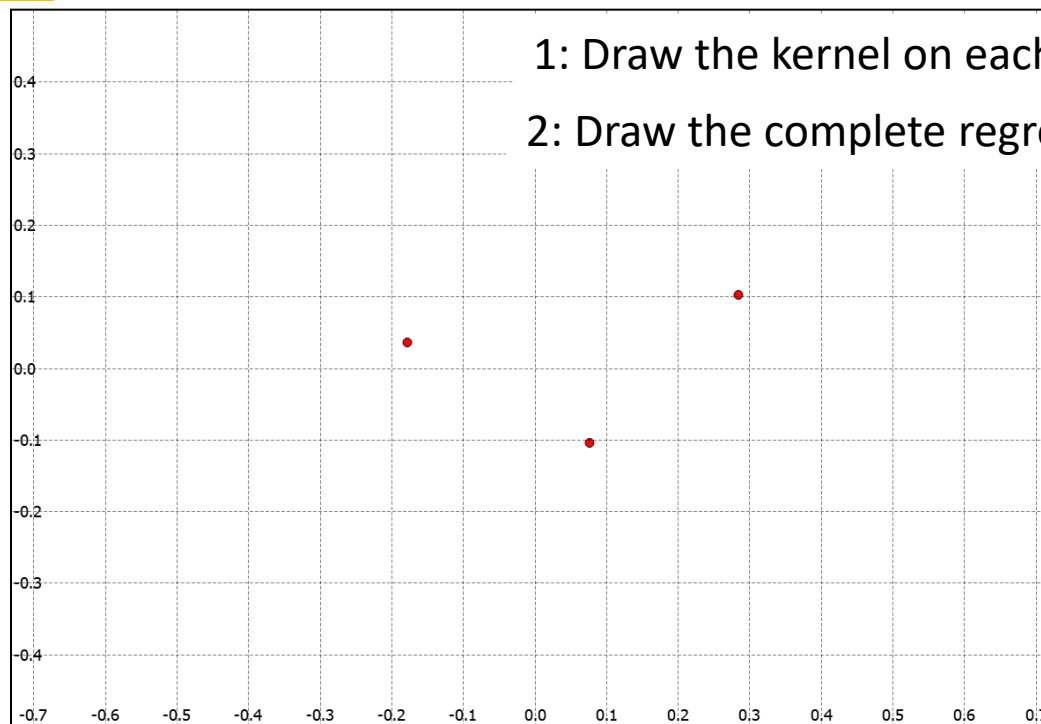
Locally weighted regression

Introduce local solution:

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

$$\beta_i(x) = e^{-\|x^i - x\|^2}$$

Also closed-form solution, but local regression



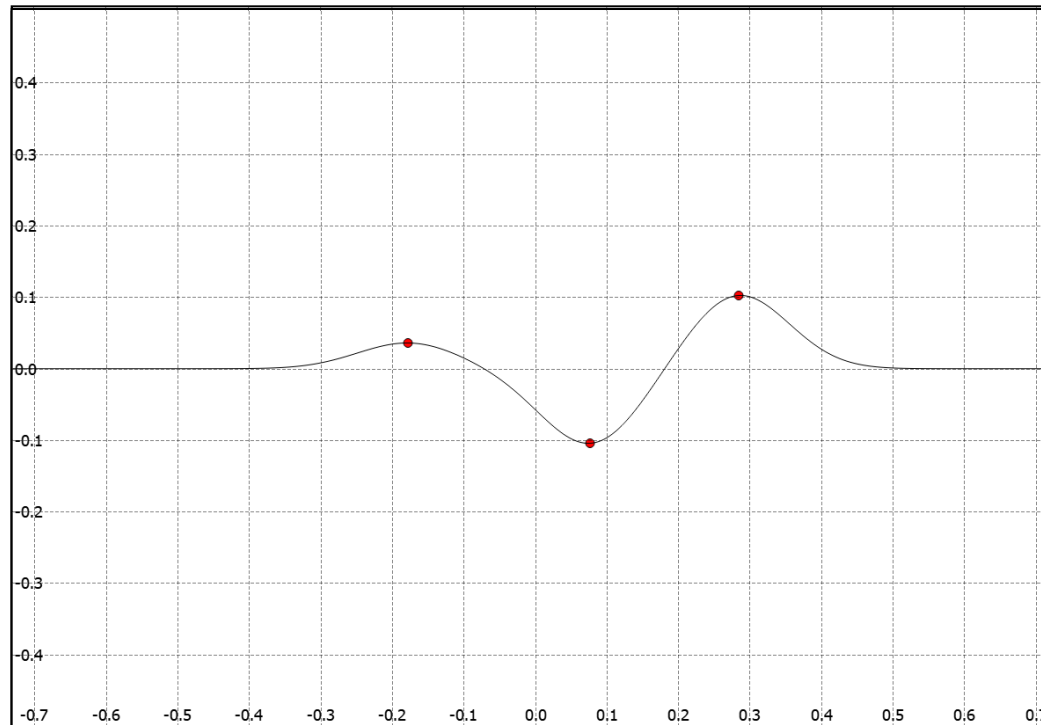
Locally weighted regression

Introduce local solution:

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

$$\beta_i(x) = e^{-\|x^i - x\|^2}$$

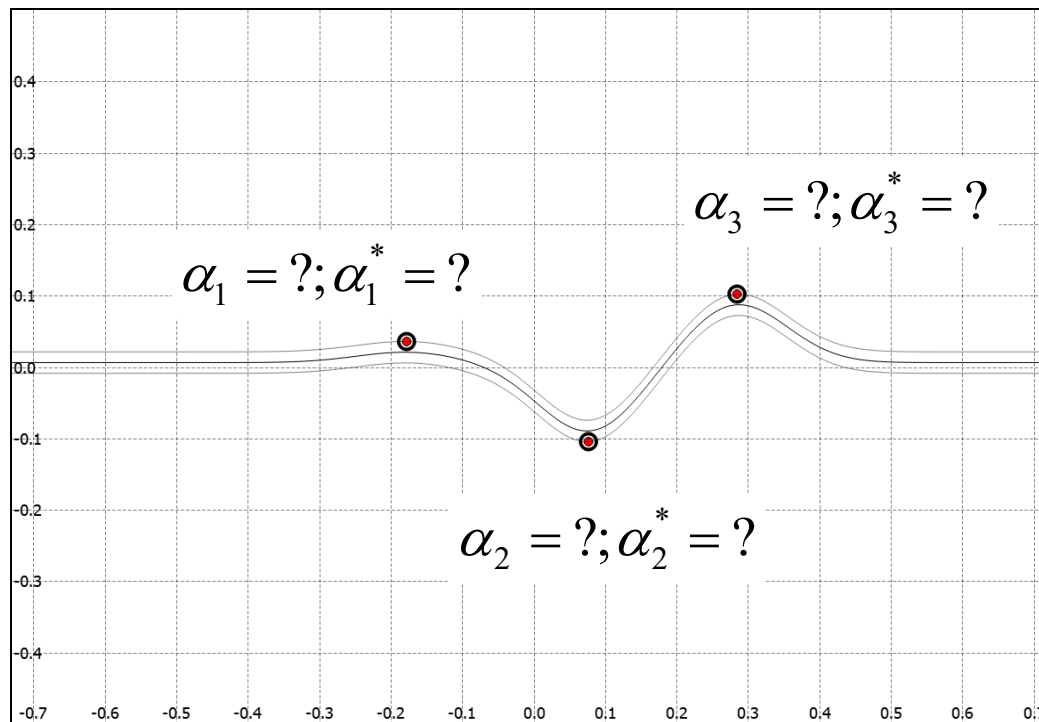
Also closed-form solution, but local regression



Support Vector Regression

SVR determines automatically which point matters for building the regression.

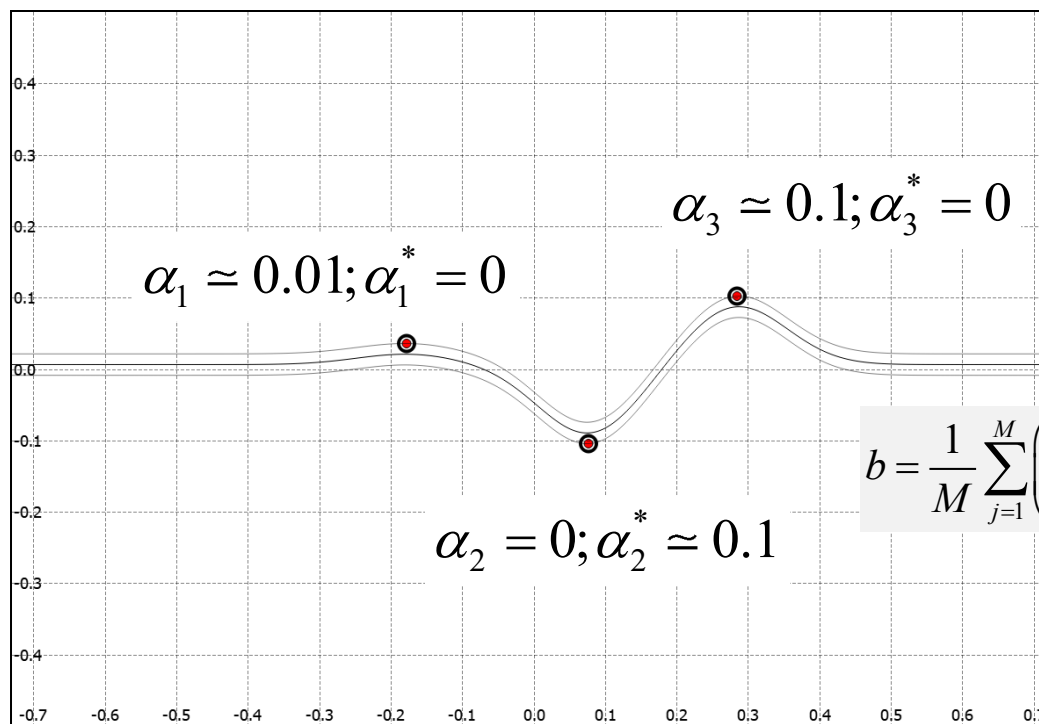
$$y = f(x) = \sum_{i=1}^M (\alpha_i - \alpha_i^*) k(x^i, x) + b$$



Support Vector Regression

Determines automatically which point matters for building the regression.

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

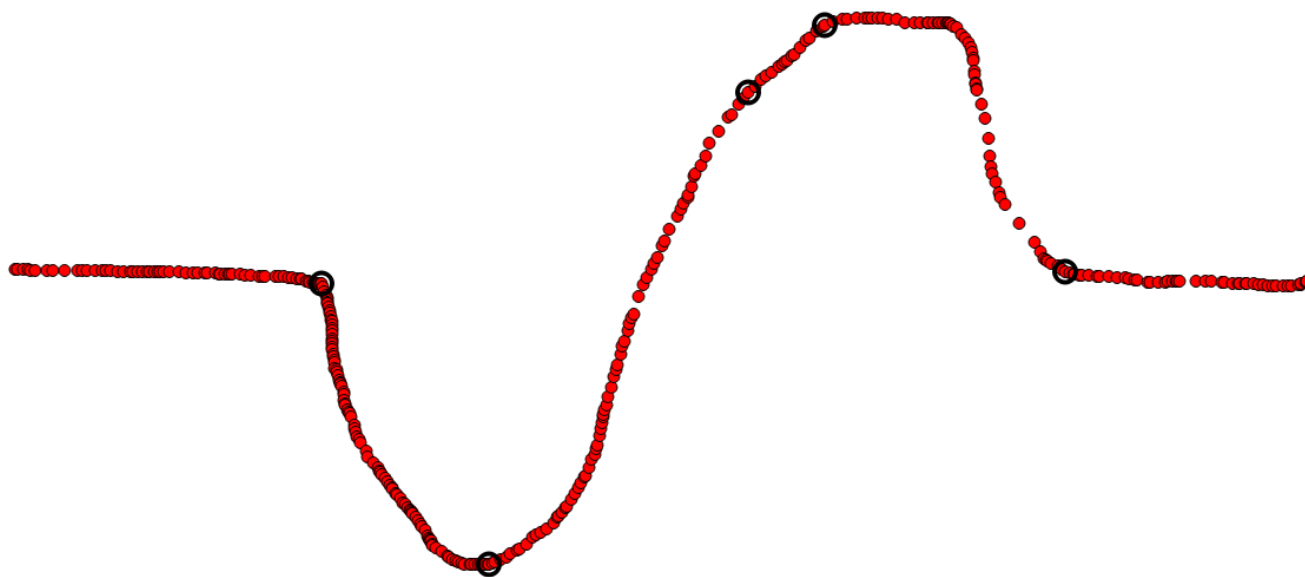


$$b \approx 0$$

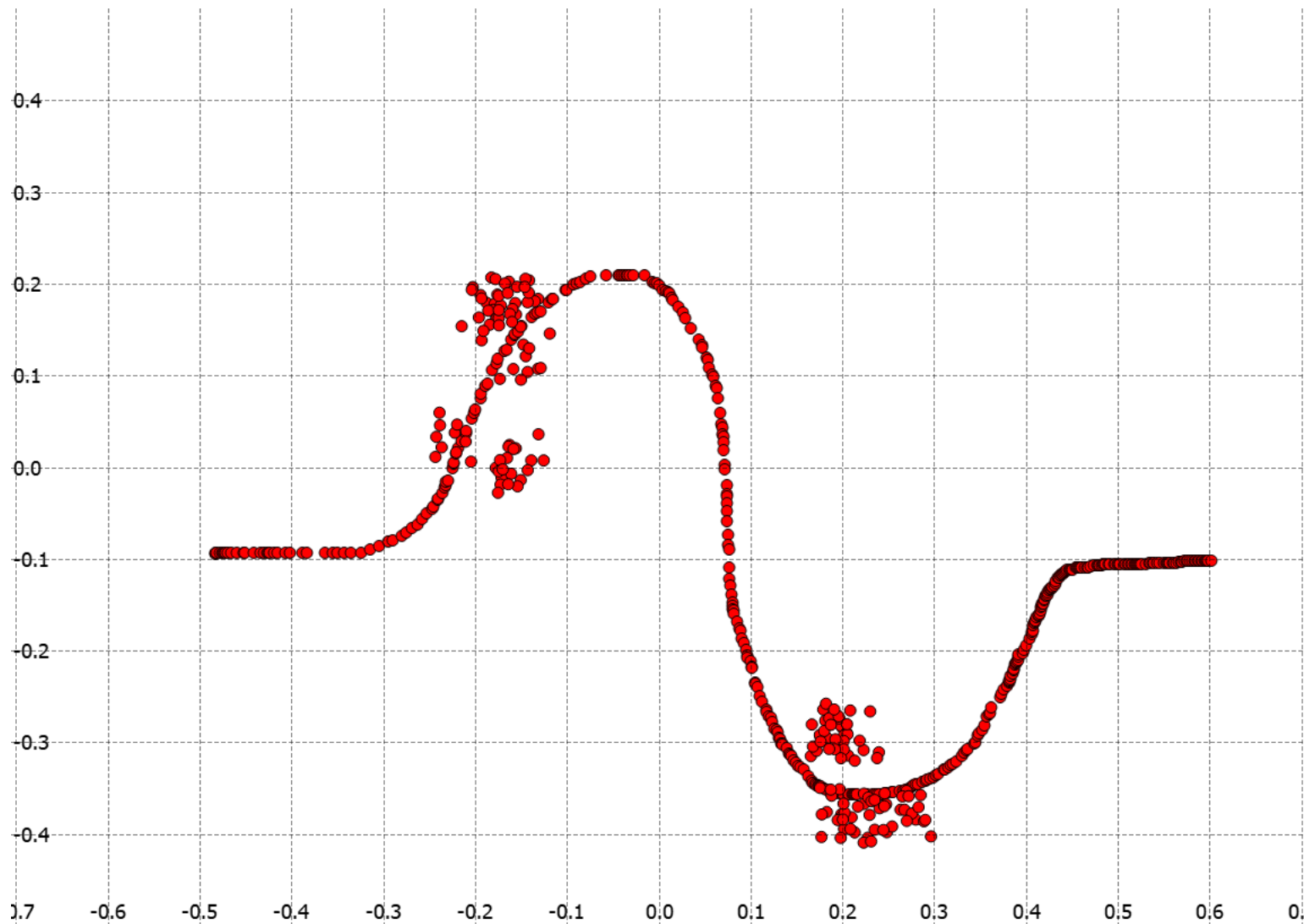
$$b = \frac{1}{M} \sum_{j=1}^M \left(y^j - \sum_{i=1}^M (\alpha_i - \alpha_i^*) k(x^j, x^i) \right)$$

Support Vector Regression

SVR determines the weight associated to each point. If it sets a weight of zero, it discards the point. This allows to have a sparse representation of the curve.



SVR : noise sensitivity



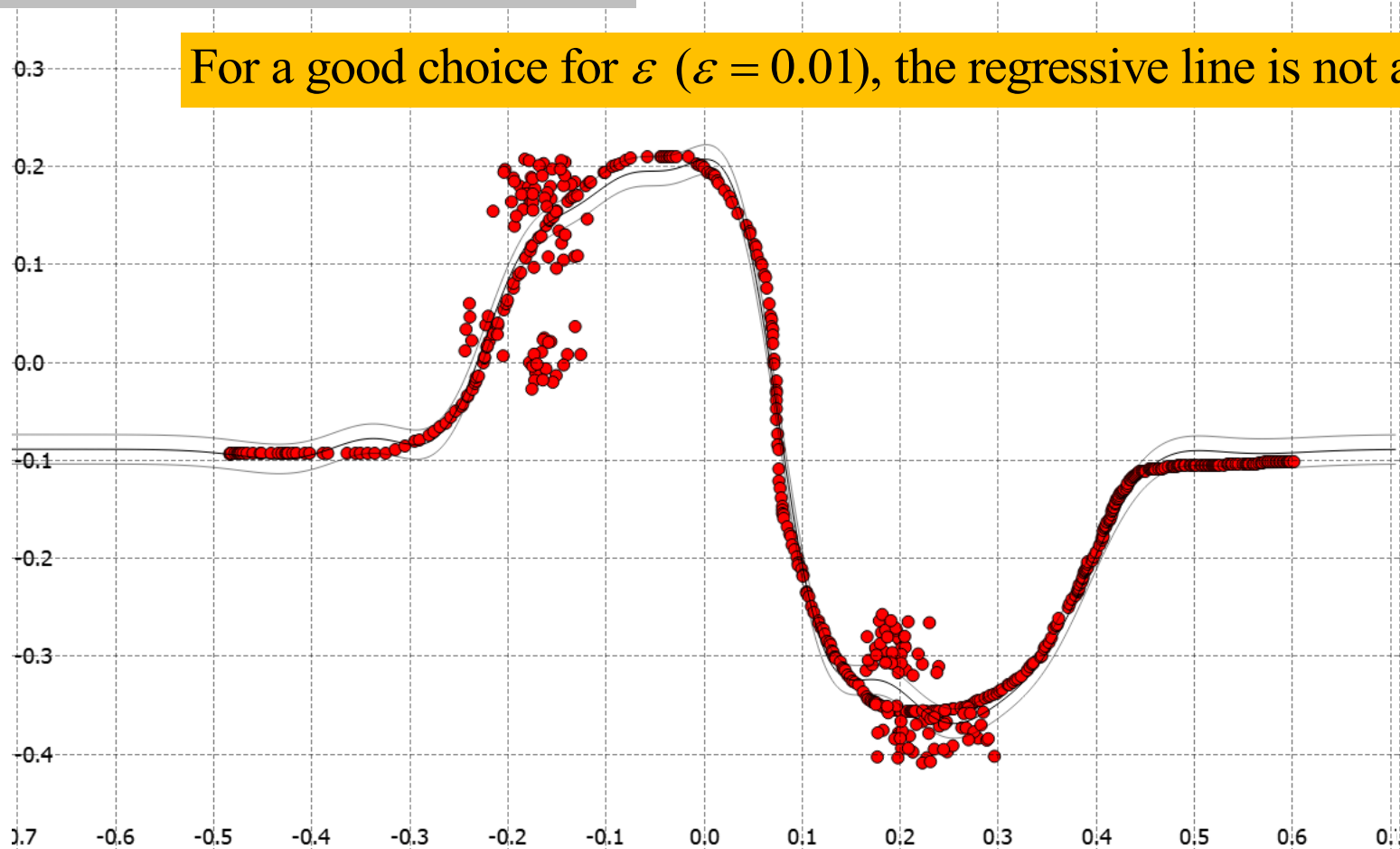
How would SVR handle this noise?

SVR : noise sensitivity

SVR looks for a function f , s.t.

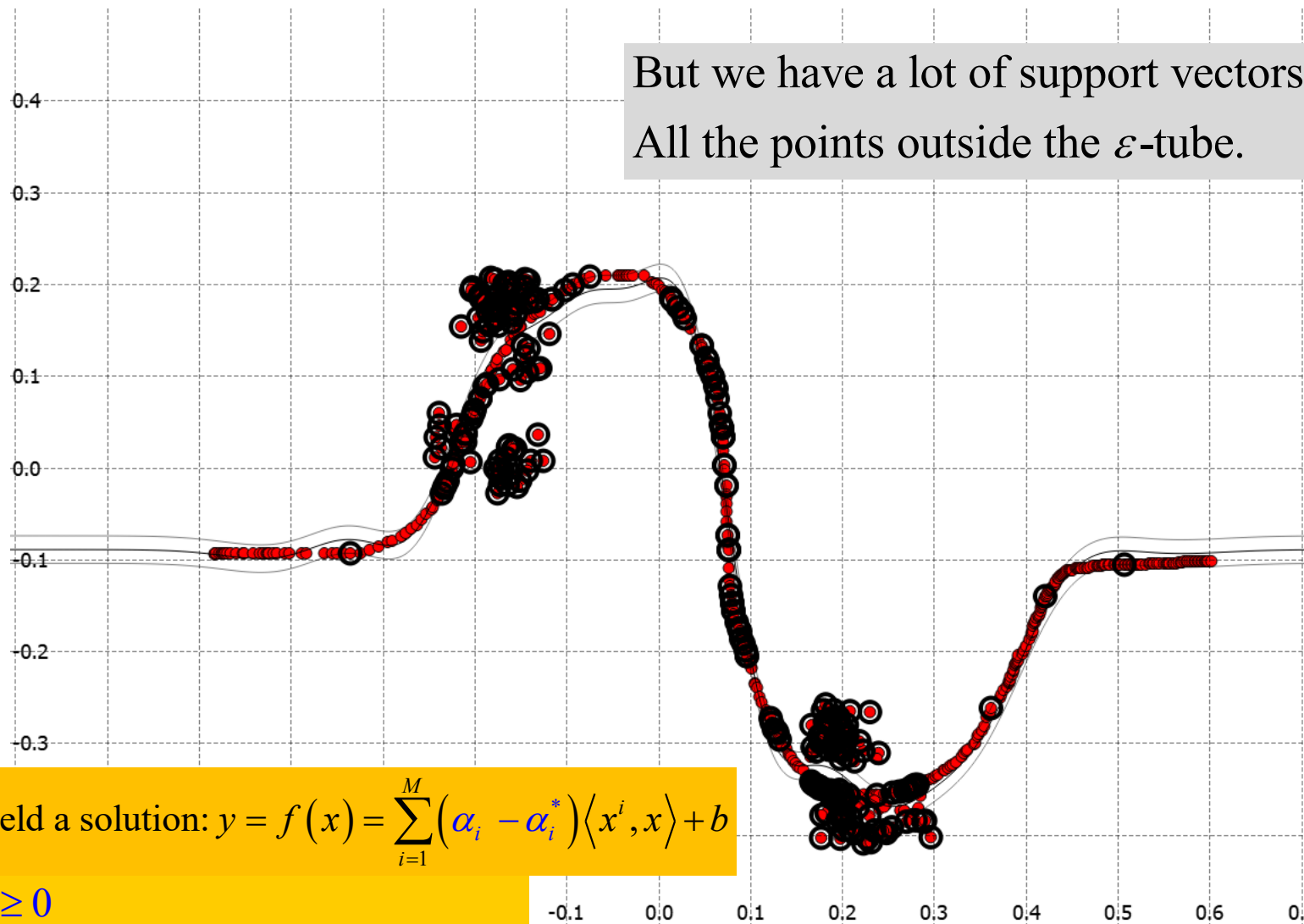
$$f(x) - y \leq \varepsilon, \quad \varepsilon \geq 0: \text{noise estimate}$$

For a good choice for ε ($\varepsilon = 0.01$), the regressive line is not affected.



SVR : noise sensitivity

But we have a lot of support vectors.
All the points outside the ε -tube.



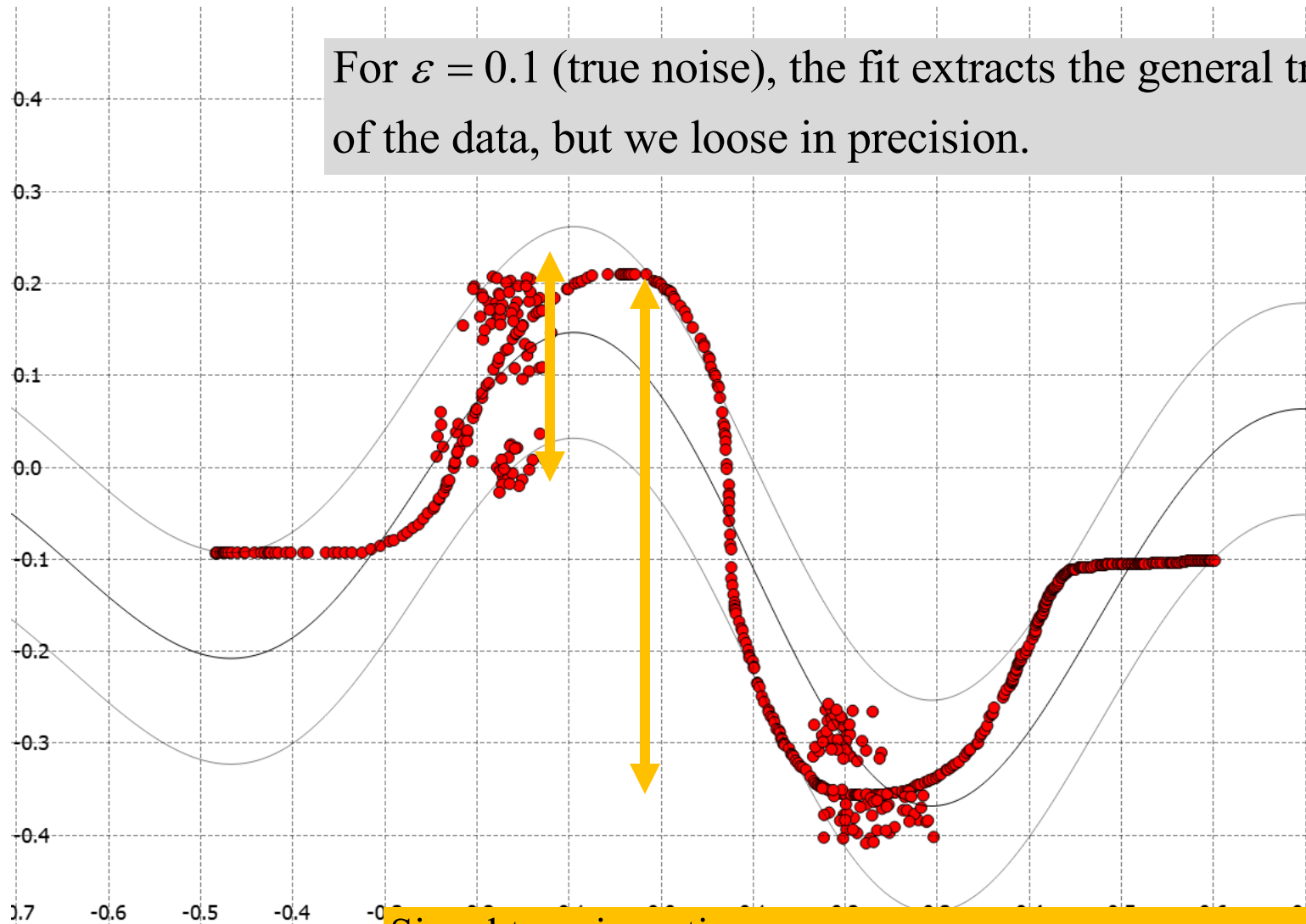
SVR yield a solution: $y = f(x) = \sum_{i=1}^M (\alpha_i - \alpha_i^*) \langle x^i, x \rangle + b$

$$\alpha_i, \alpha_i^* \geq 0$$

α_i or $\alpha_i^* > 0$ for points outside ε -tube

Regression: noise

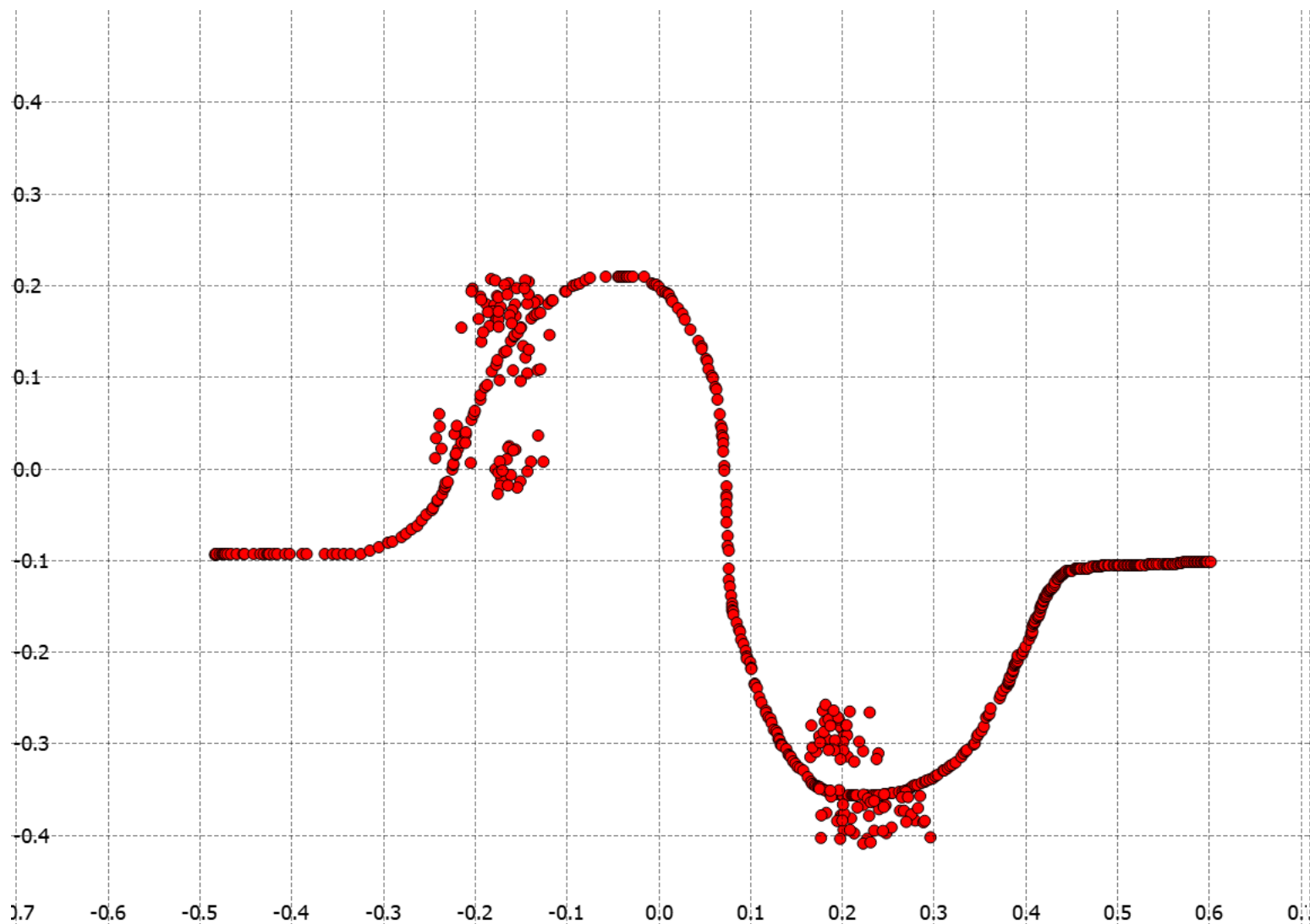
For $\varepsilon = 0.1$ (true noise), the fit extracts the general trend of the data, but we loose in precision.



Signal to noise ratio :

variance of y should be much larger than noise magnitude

Regression: noise



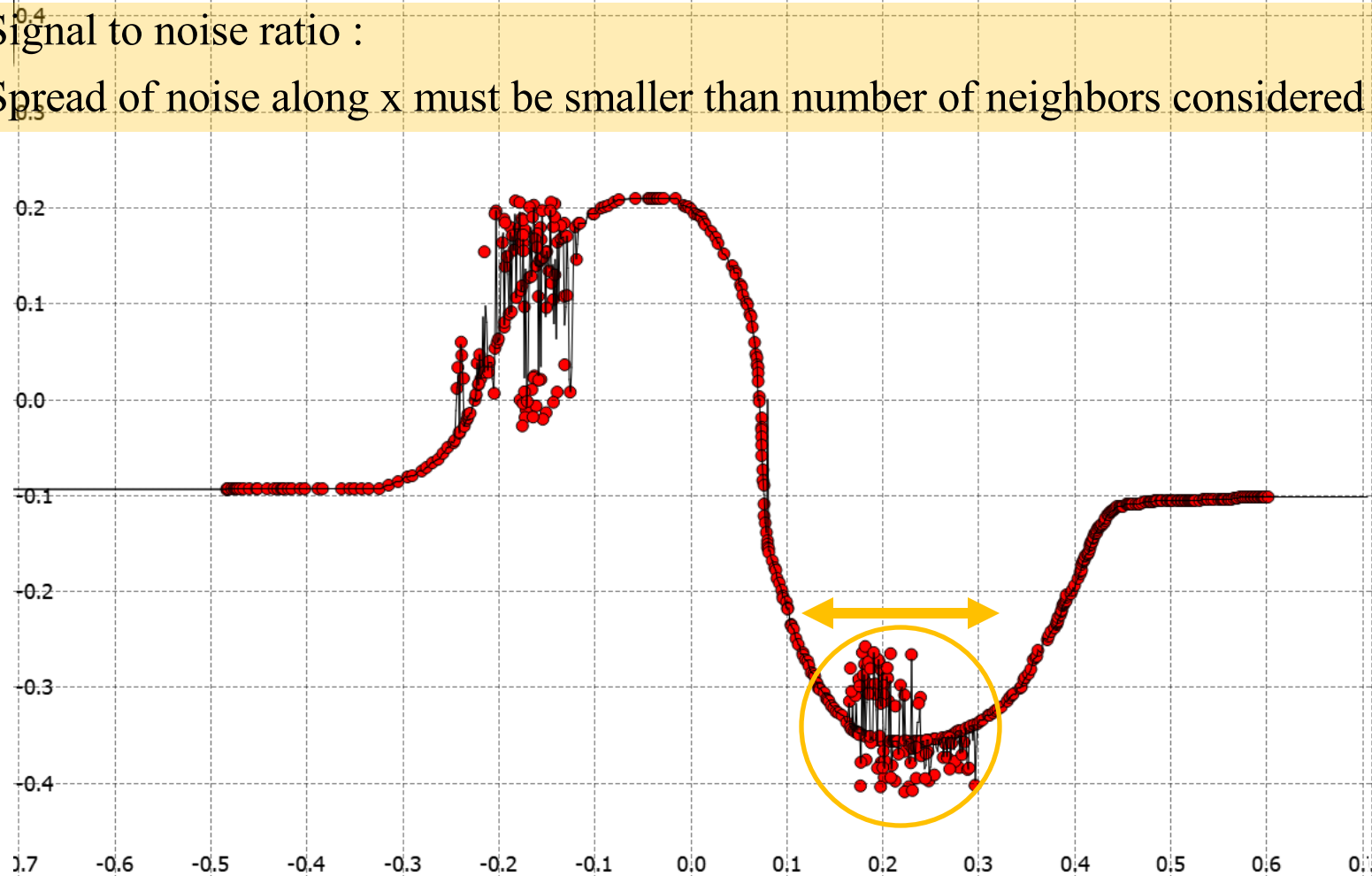
How would KNN handles this noise?

Regression: noise

Noise leads to strong fluctuation for a small K , here $K=1$

Signal to noise ratio :

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

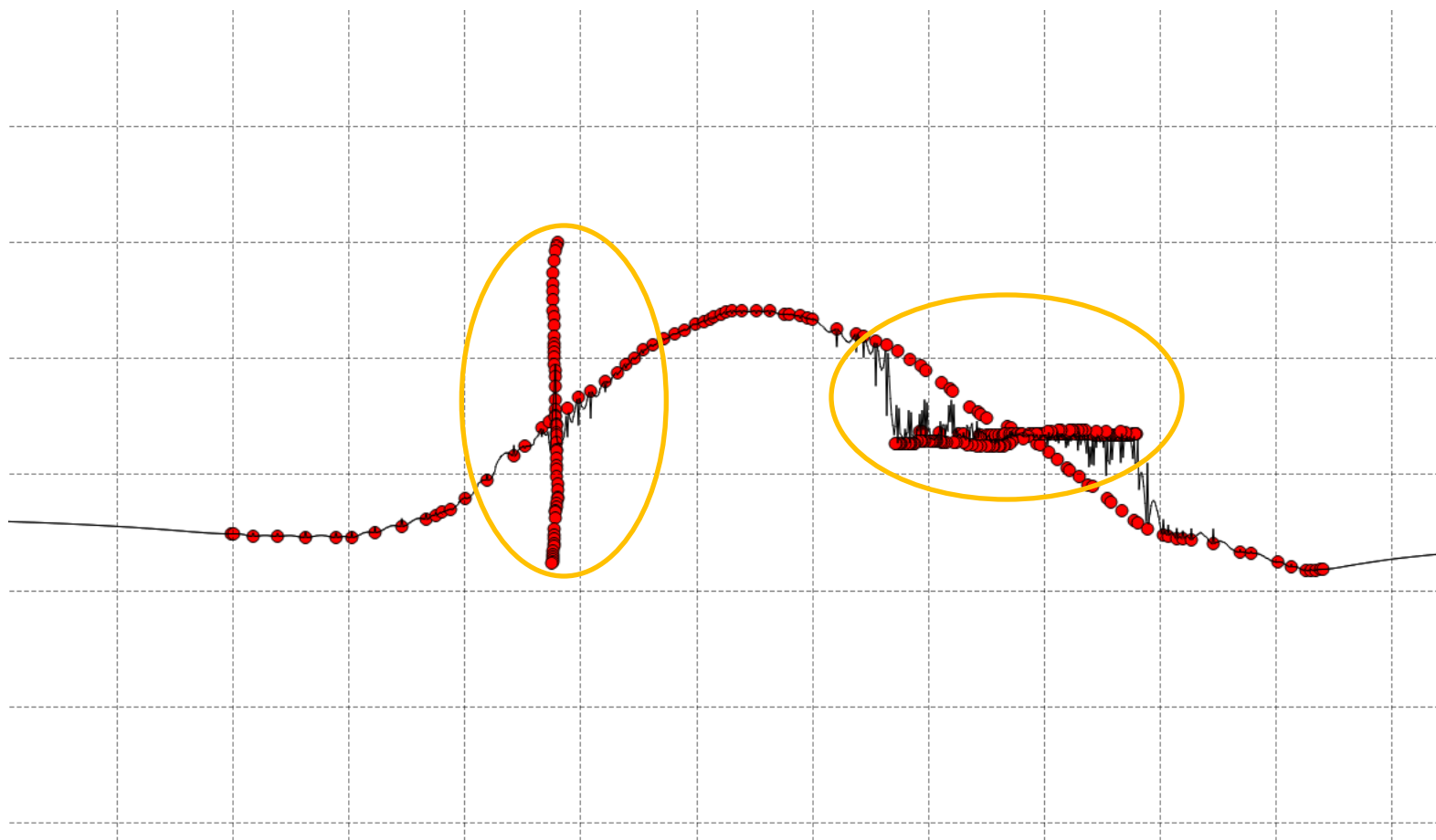


How would KNN handles this noise?

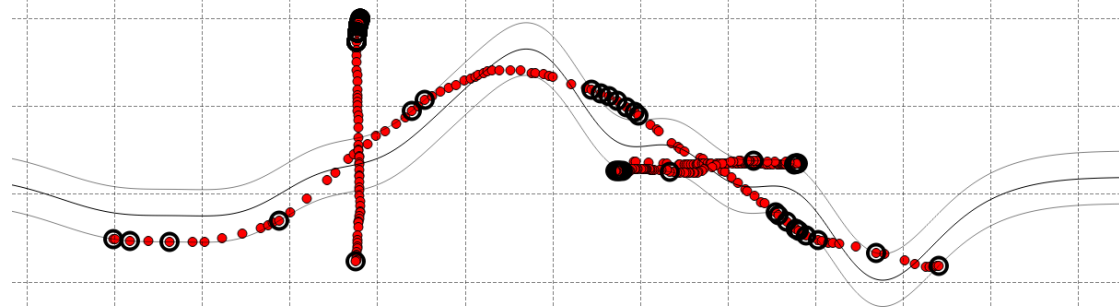
Regression: noise

Signal to noise ratio :

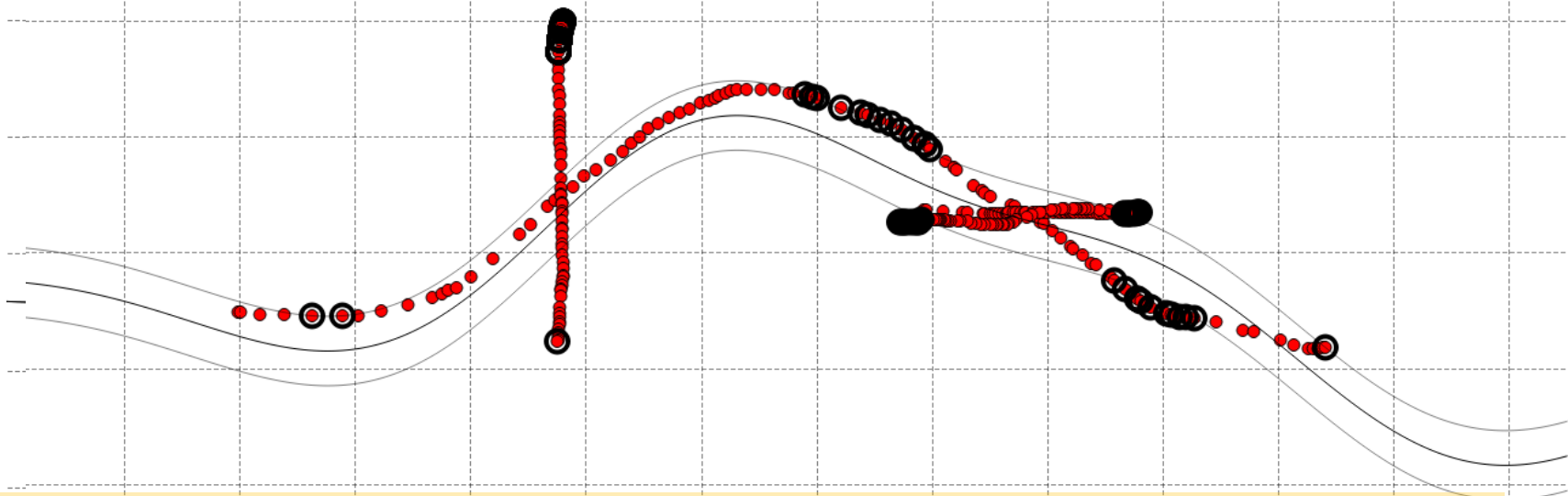
Spread of noise along x must be smaller than number of neighbors considered (K)



Reg

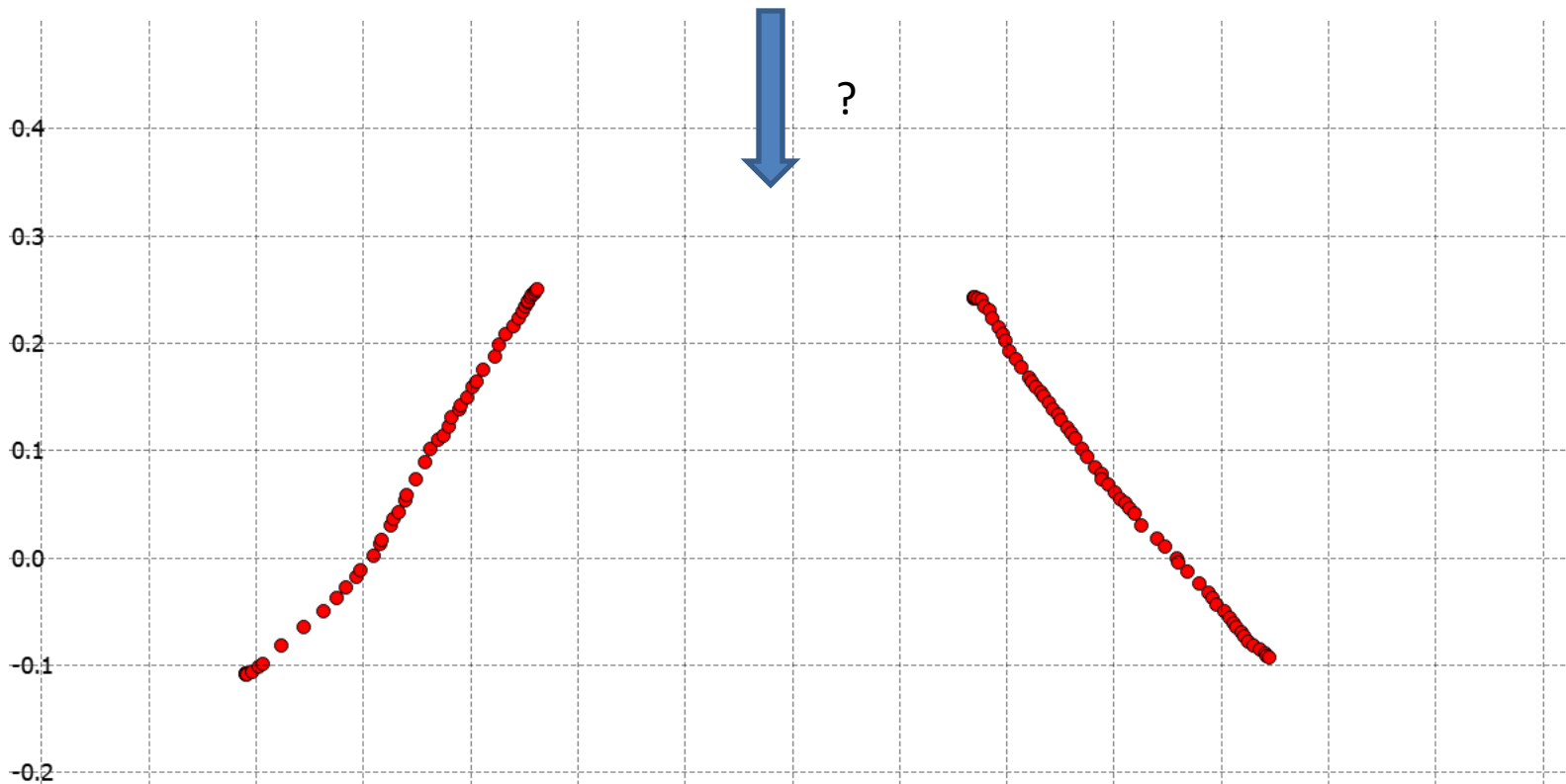
Very small σ .

The smaller the kernel width σ , the larger the influence of the noise on x .



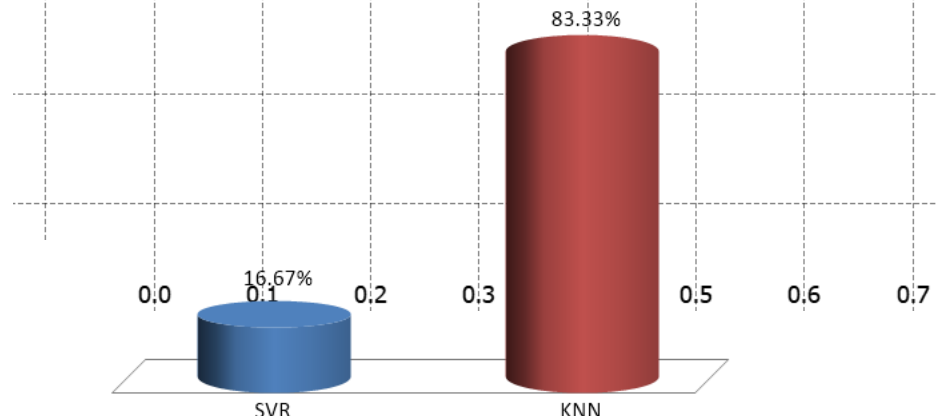
With SVR, the effect of noise on x and y is influenced both by ε and the kernel width σ .

Regression: interpolation

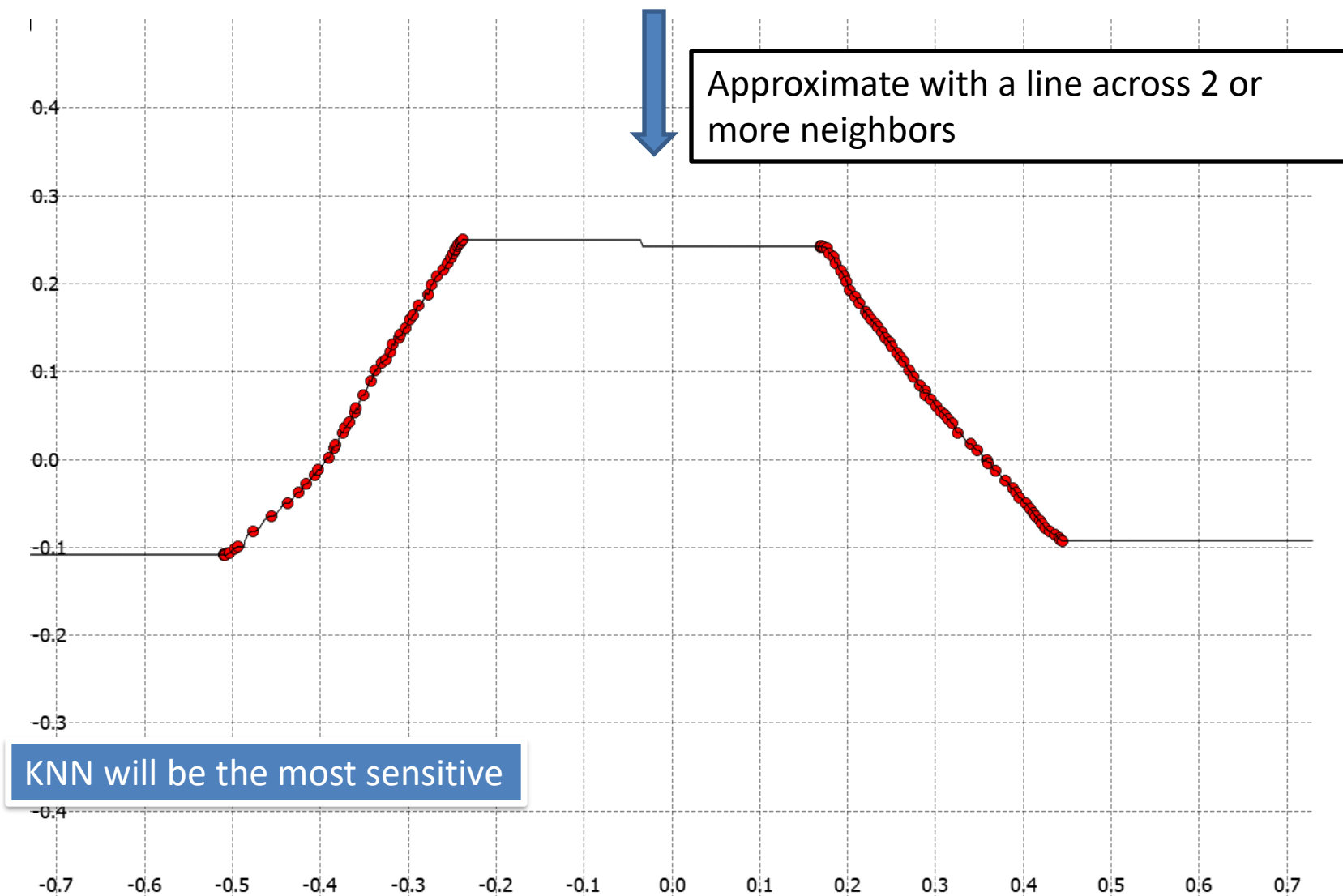


Which technique will be most sensitive to missing data?

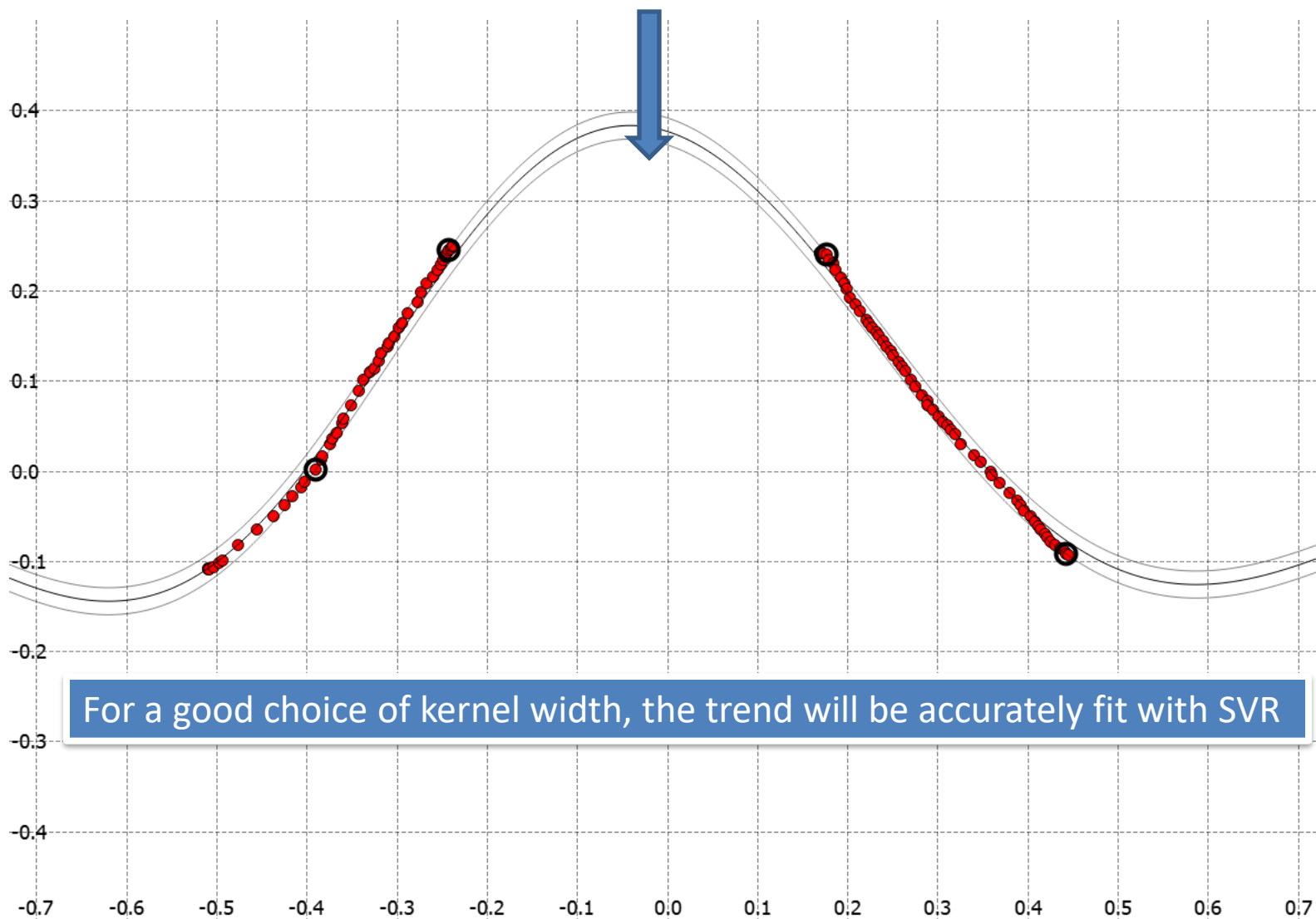
- A. SVR
- B. KNN



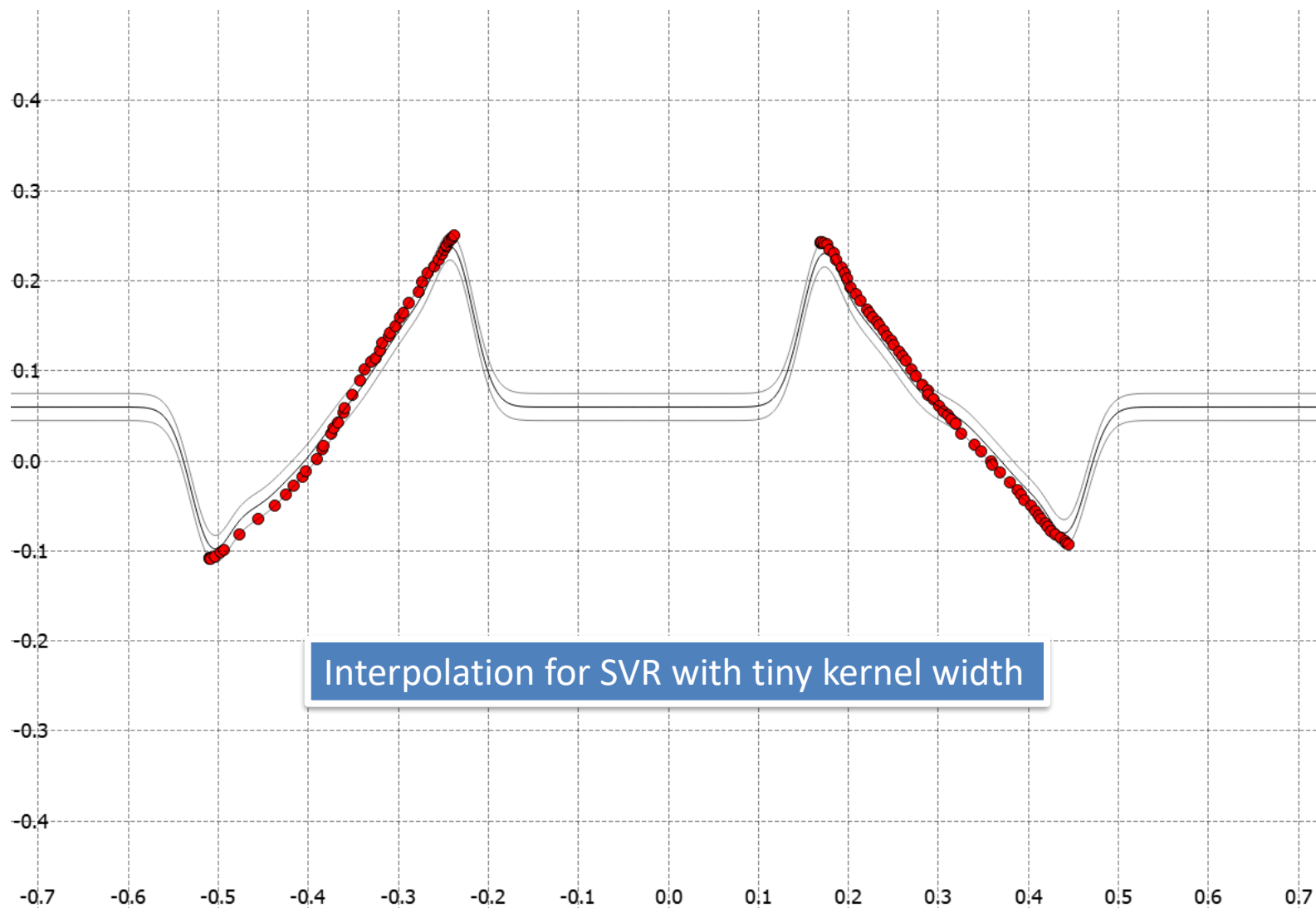
Regression: interpolation



Regression: interpolation



Regression: interpolation




Gaussian mixture regression

We must first learn the joint distribution

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


$$\alpha_k \in [0, 1]$$

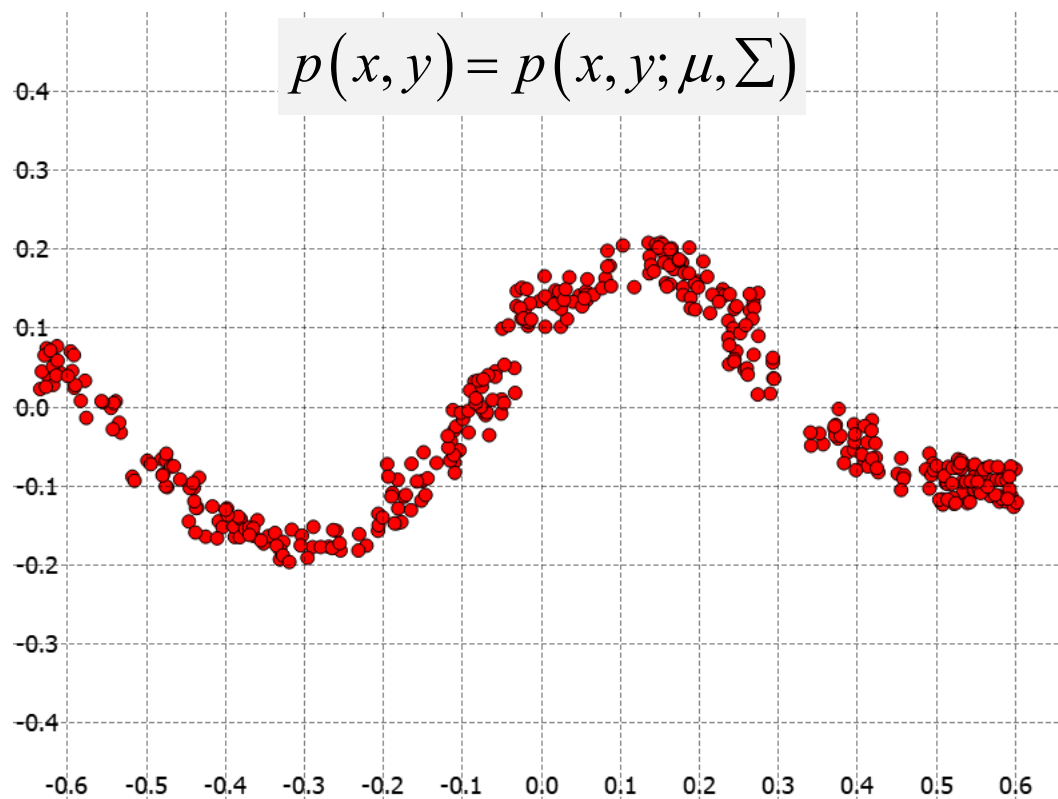

$$p(x, y; \mu^k, \Sigma^k) = N(\mu^k, \Sigma^k)$$

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

And then we compute the regressive signal:

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

Gaussian mixture regression: 1 Gauss model



Gaussian mixture regression: 1 Gauss model

$$p(x, y) = p(x, y; \mu, \Sigma)$$

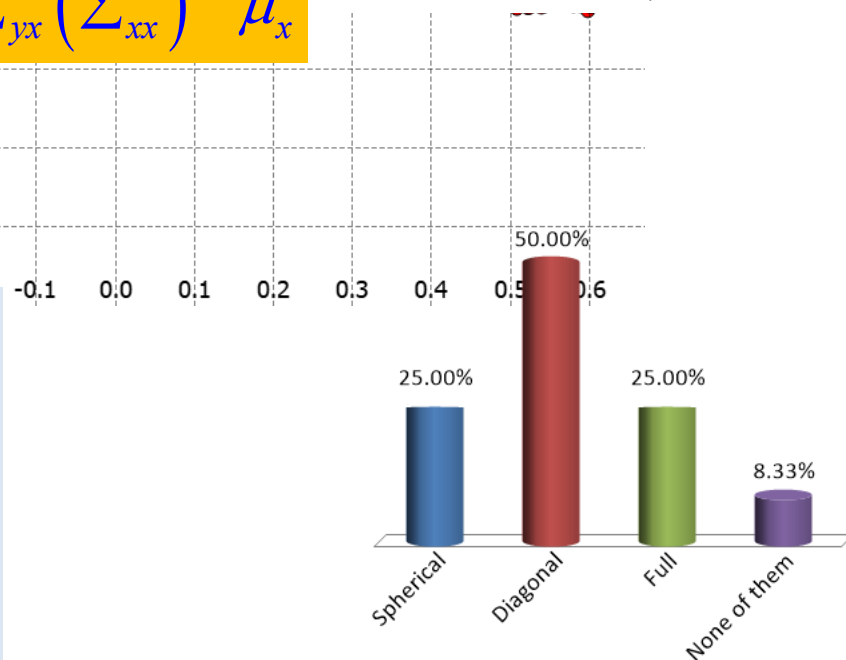
$$y = E \{ p(y | x) \} = \sum_{k=1}^K \beta_k(x) \left(\mu_y^k + \Sigma_{yx}^k \left(\Sigma_{xx}^k \right)^{-1} (x - \mu_x^k) \right)$$

$$K = 1 \Rightarrow y = \mu_y + \Sigma_{yx} \left(\Sigma_{xx} \right)^{-1} (x - \mu_x) \quad \text{with} \quad \beta_k(x) = \frac{\alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}{\sum_{k=1}^K \alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}$$

$$y = \Sigma_{yx} \left(\Sigma_{xx} \right)^{-1} x + \mu_y - \Sigma_{yx} \left(\Sigma_{xx} \right)^{-1} \mu_x$$

For which model is the regressive curve a straight line?

- A. Spherical ✓
- B. Diagonal ✓
- C. Full ✓
- D. None of them



Gaussian mixture regression: 1 Gauss model

$$p(x, y) = p(x, y; \mu, \Sigma)$$

We compute the regressive line by taking the expectation over the conditional

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

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

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

$$\Leftrightarrow y = \sum_{yx} \left(\sum_{xx} \right)^{-1} x + \mu_y - \sum_{yx} \left(\sum_{xx} \right)^{-1} \mu_x$$

Linear equation $y = Ax + b$

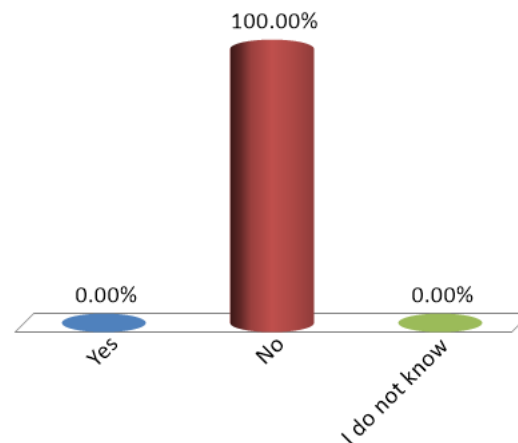
$$A = \sum_{yx} \left(\sum_{xx} \right)^{-1},$$

$$b = \mu_y - \sum_{yx} \left(\sum_{xx} \right)^{-1} \mu_x$$

Gaussian mixture regression: 1 Gauss model

Will the solutions differ when using spherical or diagonal covariance matrices?

- A. Yes
- B. No ✓
- C. I do not know



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

$$\Leftrightarrow y = \sum_{yx} \left(\sum_{xx} \right)^{-1} x + \mu_y - \sum_{yx} \left(\sum_{xx} \right)^{-1} \mu_x$$

Linear equation $y = Ax + b$

$$A = \sum_{yx} \left(\sum_{xx} \right)^{-1},$$

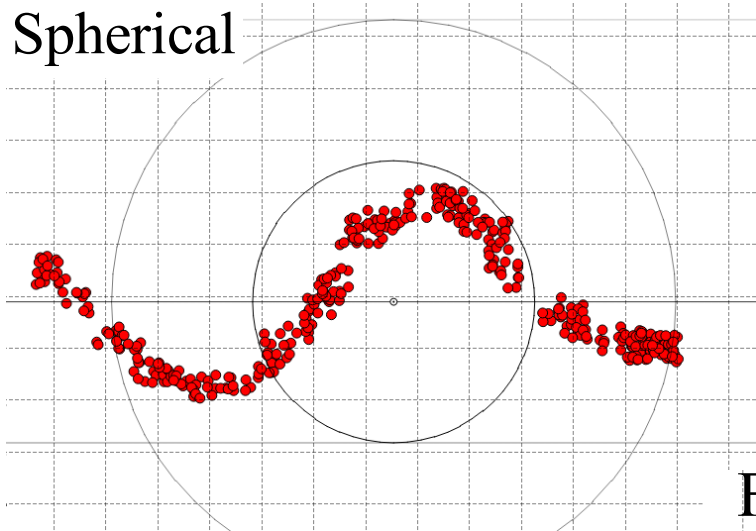
$$b = \mu_y - \sum_{yx} \left(\sum_{xx} \right)^{-1} \mu_x$$

Gaussian mixture regression: 1 Gauss model

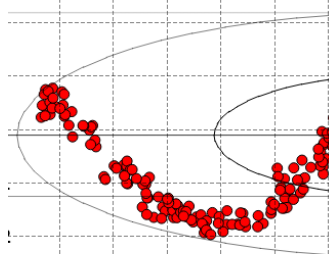
$$\Leftrightarrow y = \Sigma_{yx} \left(\Sigma_{xx} \right)^{-1} x + \mu_y - \Sigma_{yx} \left(\Sigma_{xx} \right)^{-1} \mu_x$$

$$\Sigma = \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_y \end{bmatrix}$$

Spherical



Diagonal



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

when x, y are unidimensional

Σ_x, Σ_y are scalars,

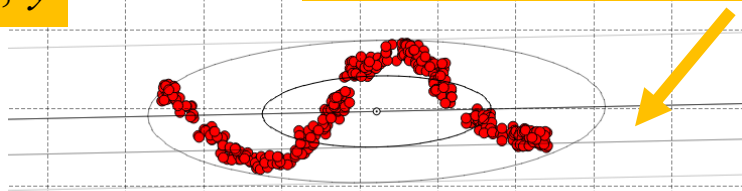
$\Sigma_x = \Sigma_y$ in the spherical case.

Full

$$\Sigma_x = \sigma_x I, \Sigma_y = \sigma_y I, \Sigma_{xy} = 0$$

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

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



$\Sigma_{xy} \neq 0$: entail correlation between x and y .

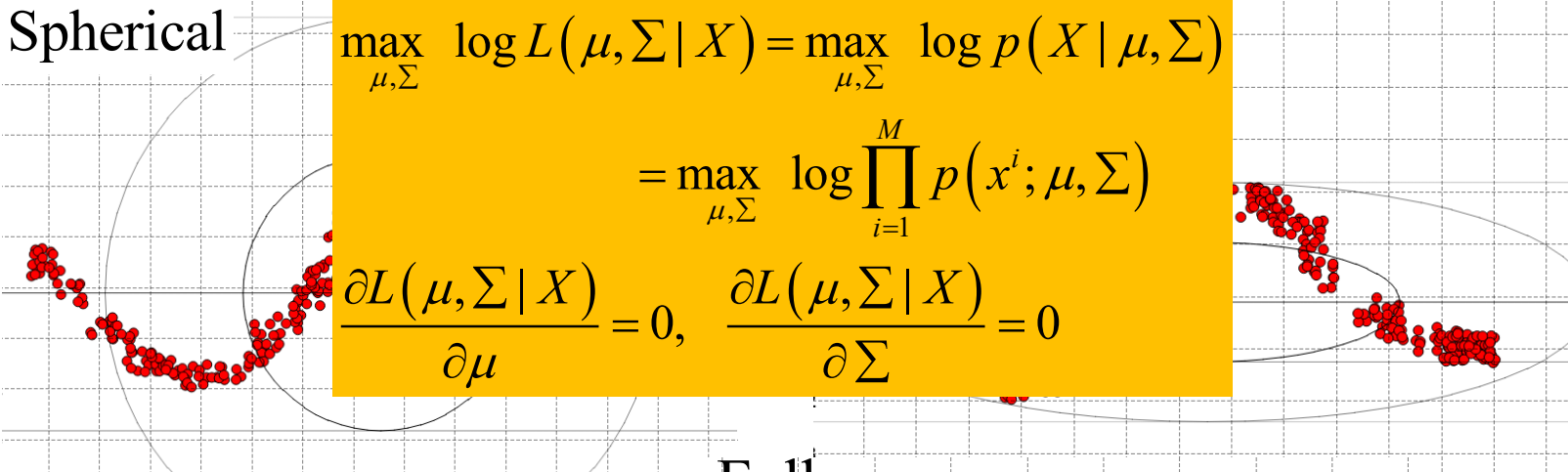
GMR – 1 Gauss fct: Uniqueness of solution

Solution of maximum likelihood problem:

$$\begin{aligned} \max_{\mu, \Sigma} \log L(\mu, \Sigma | X) &= \max_{\mu, \Sigma} \log p(X | \mu, \Sigma) \\ &= \max_{\mu, \Sigma} \log \prod_{i=1}^M p(x^i; \mu, \Sigma) \end{aligned}$$

$$\frac{\partial L(\mu, \Sigma | X)}{\partial \mu} = 0, \quad \frac{\partial L(\mu, \Sigma | X)}{\partial \Sigma} = 0$$

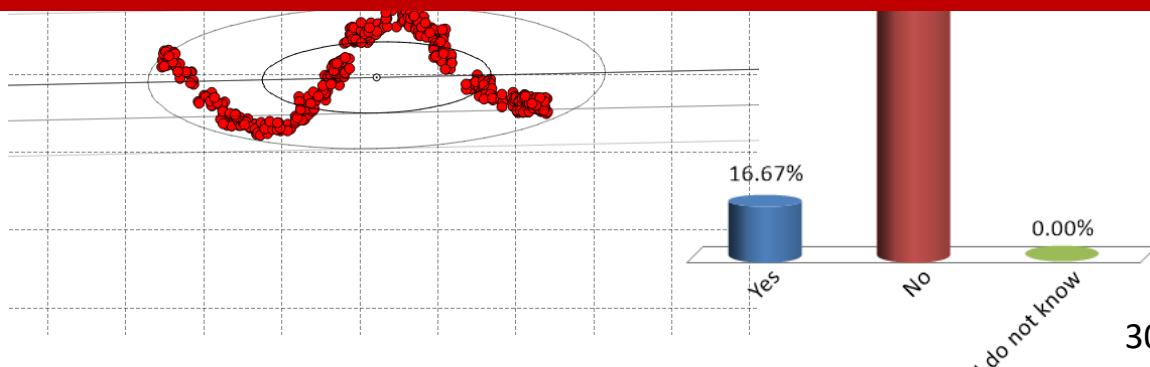
Spherical



The solution is unique and can be determined analytically:
It corresponds to the mean and variance of dataset,
see exercises of pdf lecture

Is the solution unique?

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



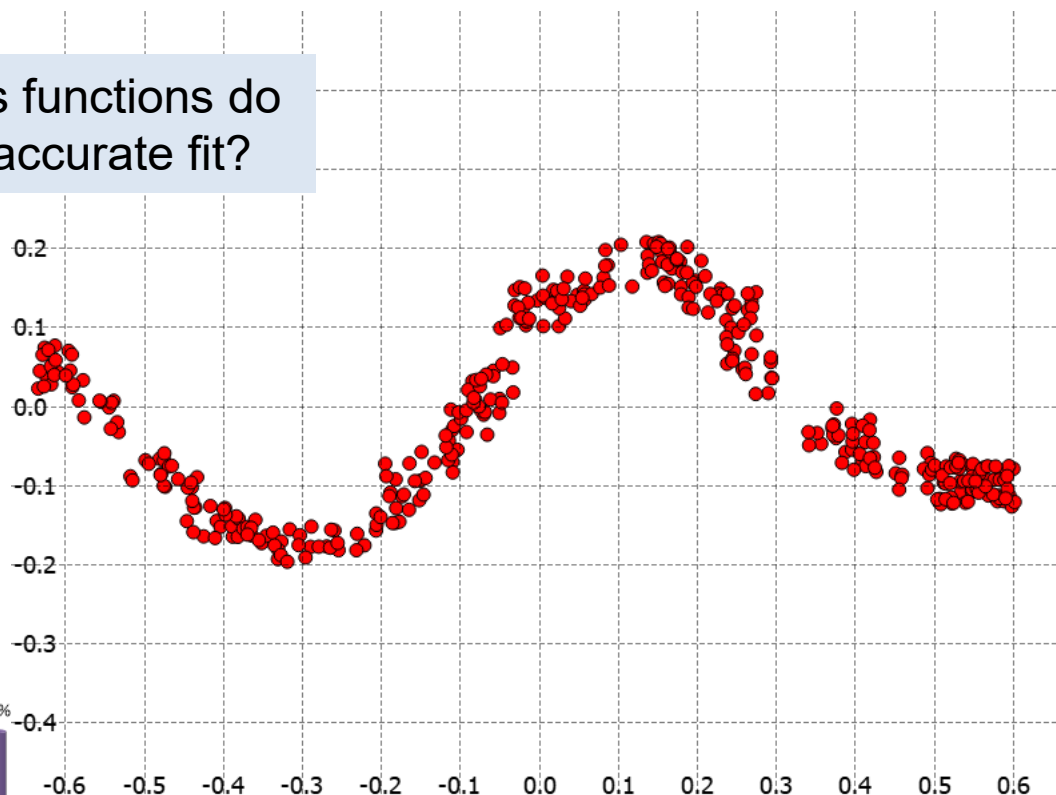
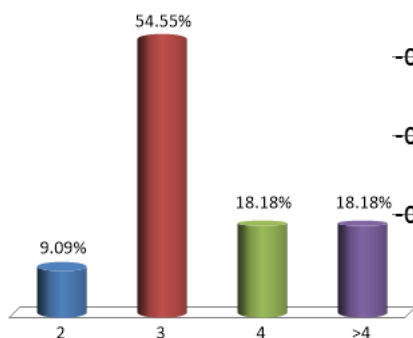
GMR: Multiple Gauss Functions

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

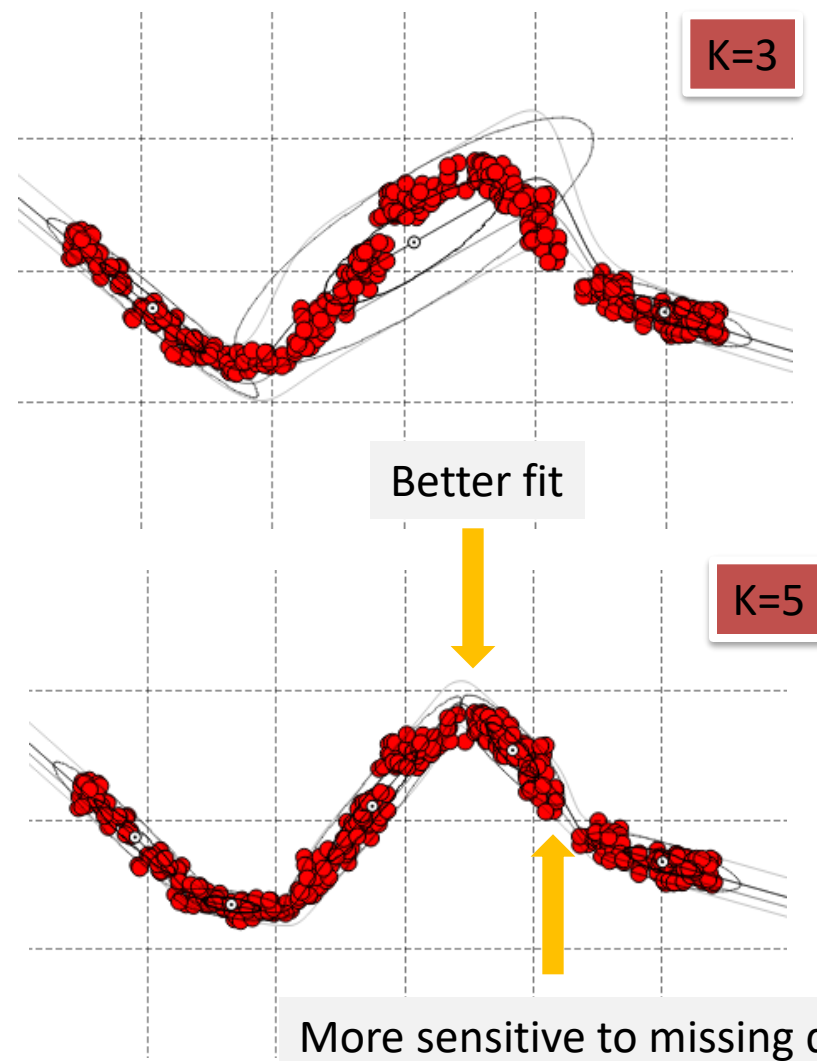
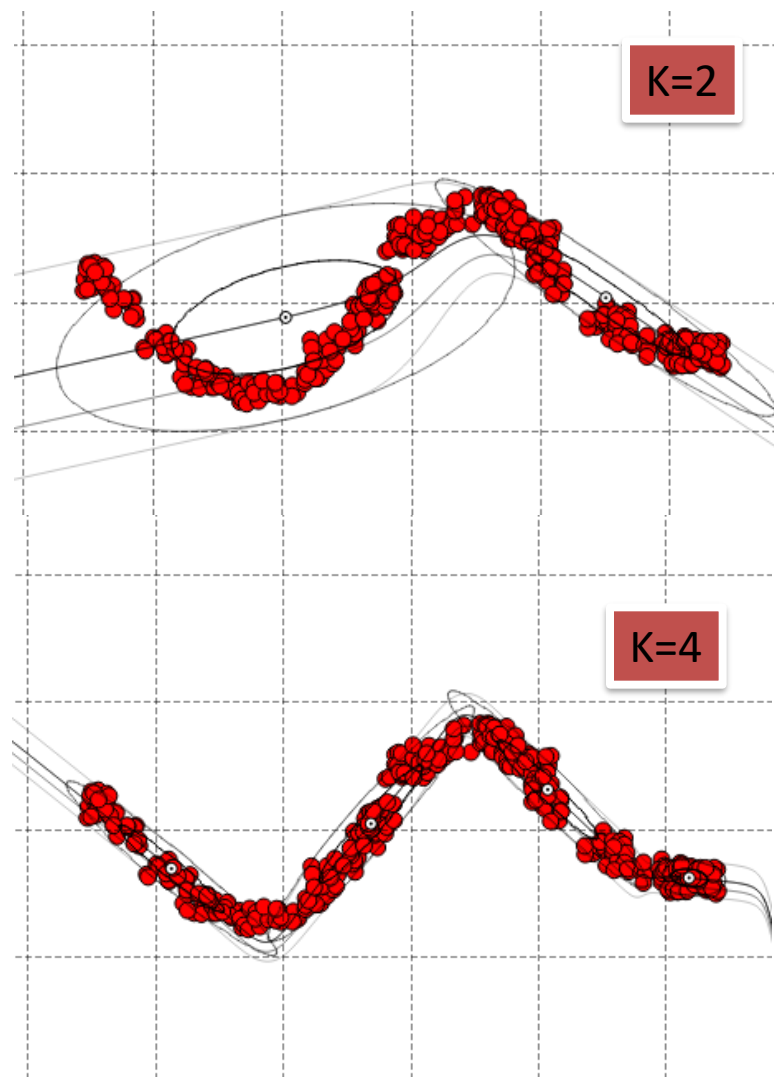
μ^i, Σ^i : mean and covariance matrix of Gaussian k.

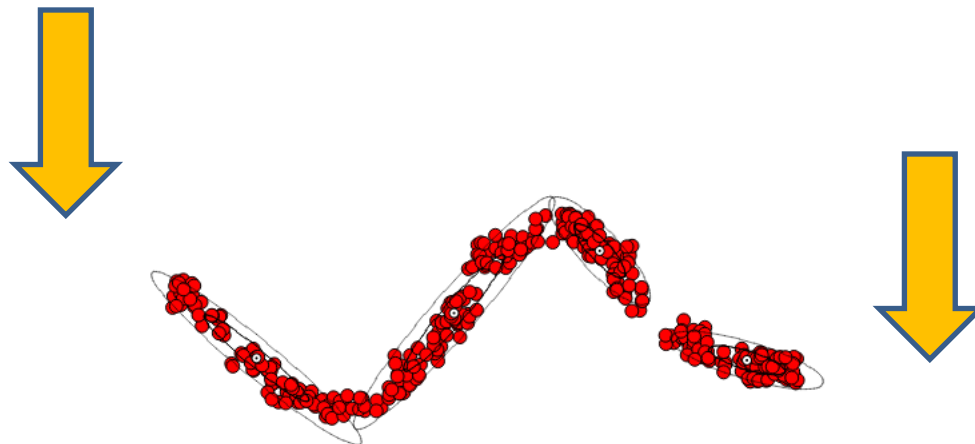
How many Gauss functions do we need for an accurate fit?

- A. 2
- B. 3
- C. 4
- D. >4



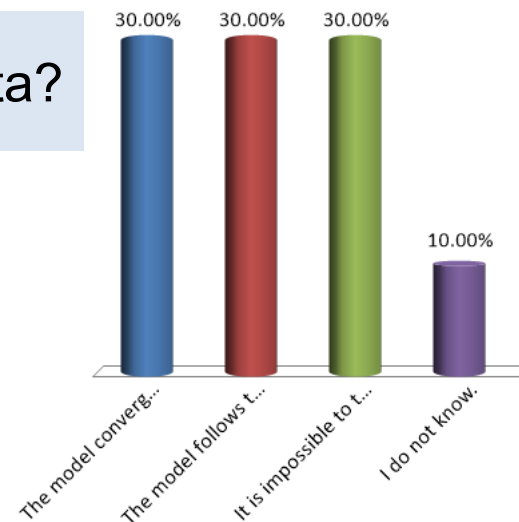
Accuracy of fit with multiple Gauss functions





What does the model predict away from the data?

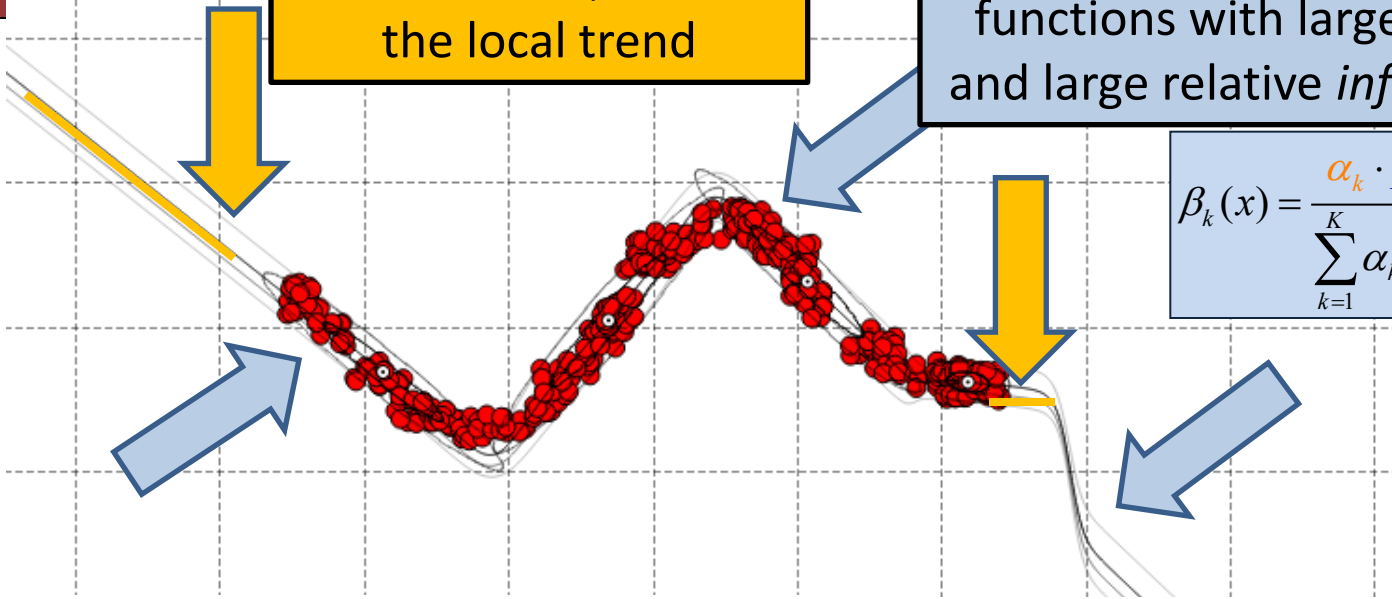
- A. The model converges to a single value.
- B. The model follows the local trend.
- C. It is impossible to tell.
- D. I do not know.



Close to data, follows the local trend

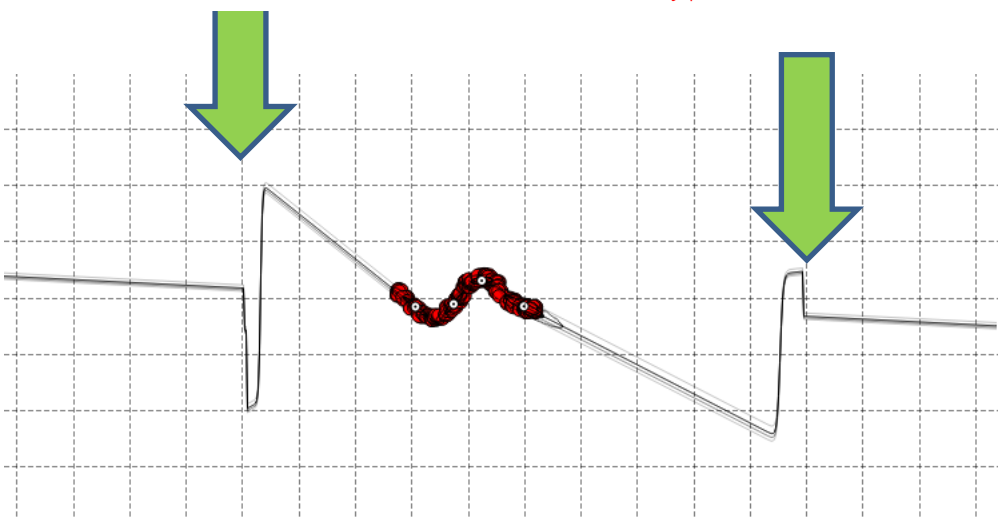
Far away, influenced by Gauss functions with large variance and large relative influence (α_k)

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



$$y = \sum_{k=1}^K \beta_k(x) \underbrace{\left(\mu_y^k + \Sigma_{yx}^k (\Sigma_{xx}^k)^{-1} (x - \mu_x^k) \right)}_{\tilde{\mu}_{y|x}^k(x)} \Rightarrow y = \sum_{k=1}^K \beta_k(x) \tilde{\mu}_{y|x}^k(x)$$

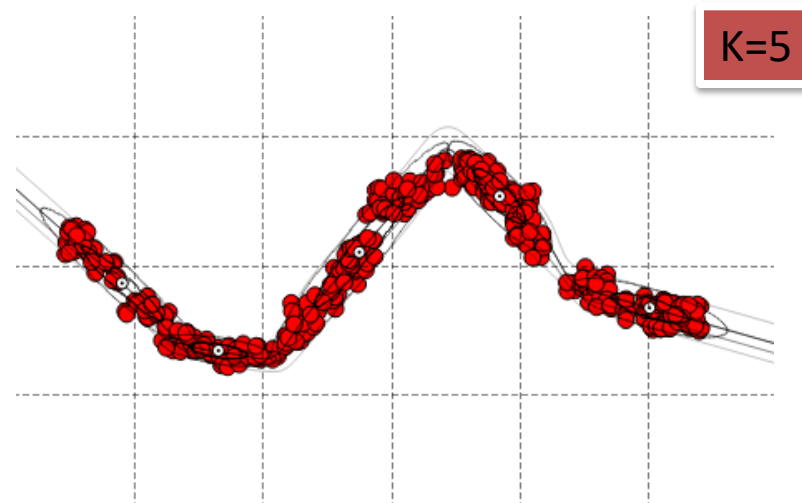
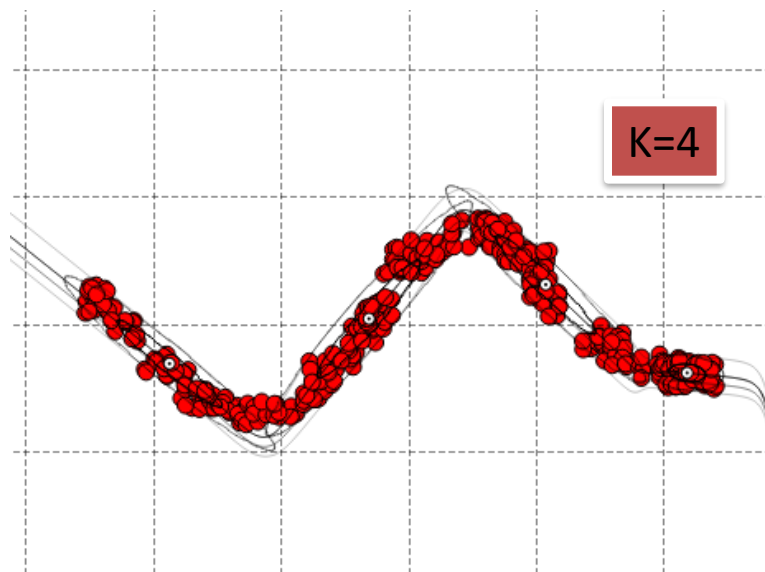
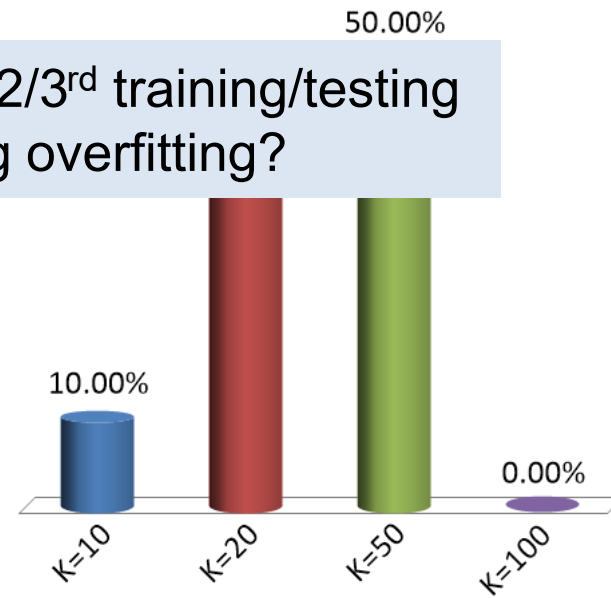
Linear combination of K local regressive models



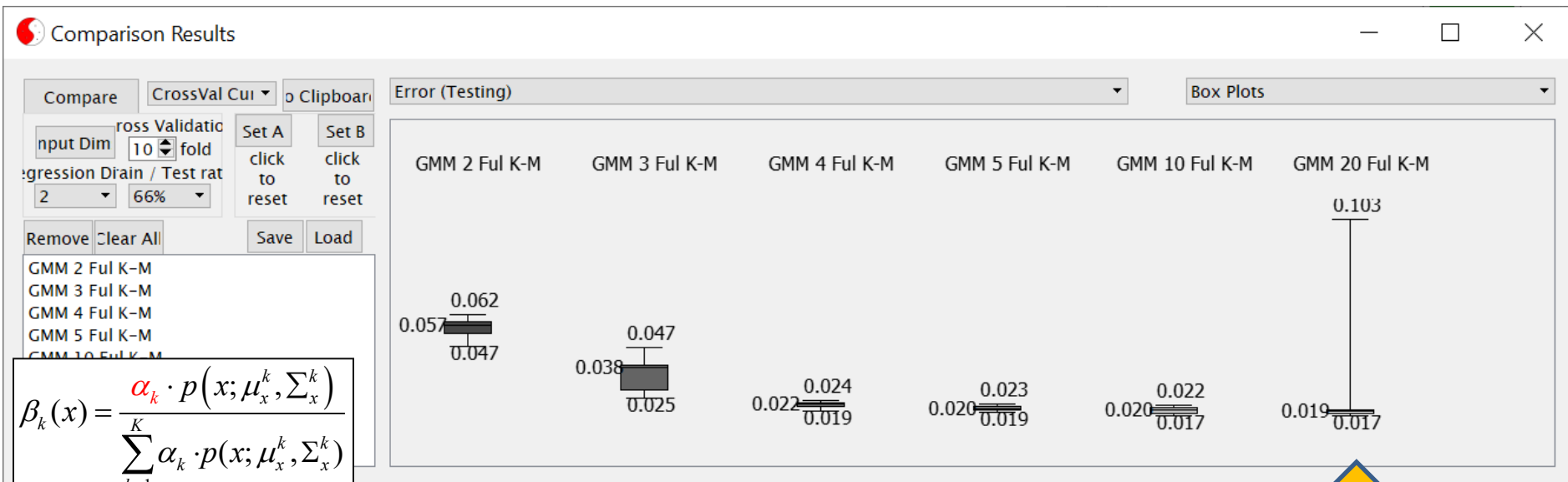
Very far from the data the prediction can be completely unexpected.

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

- A. K=10
- B. K=20
- C. K=50
- D. K=100



Overfitting with multiple Gauss functions



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

$$y = \sum_{k=1}^K \beta_k(x) \left(\mu_y^k + \sum_{yx}^k \left(\sum_{xx}^k \right)^{-1} (x - \mu_x^k) \right)$$

K Priors + $K \cdot N$ Means + $K \cdot \left(\frac{N(N+1)}{2} \right)$ Full Covariance Matrix

ratio · # datapoints · dimension = 0.66 · 100 · 2 = 132
 K=20: # parameters = 120

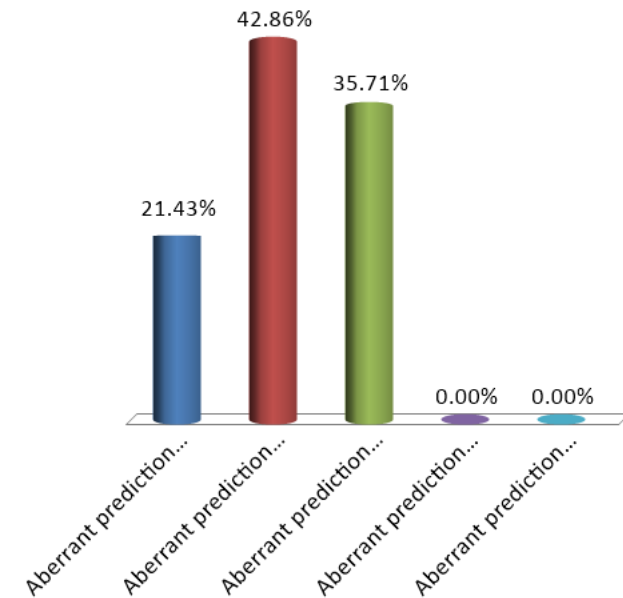


Clear overfit

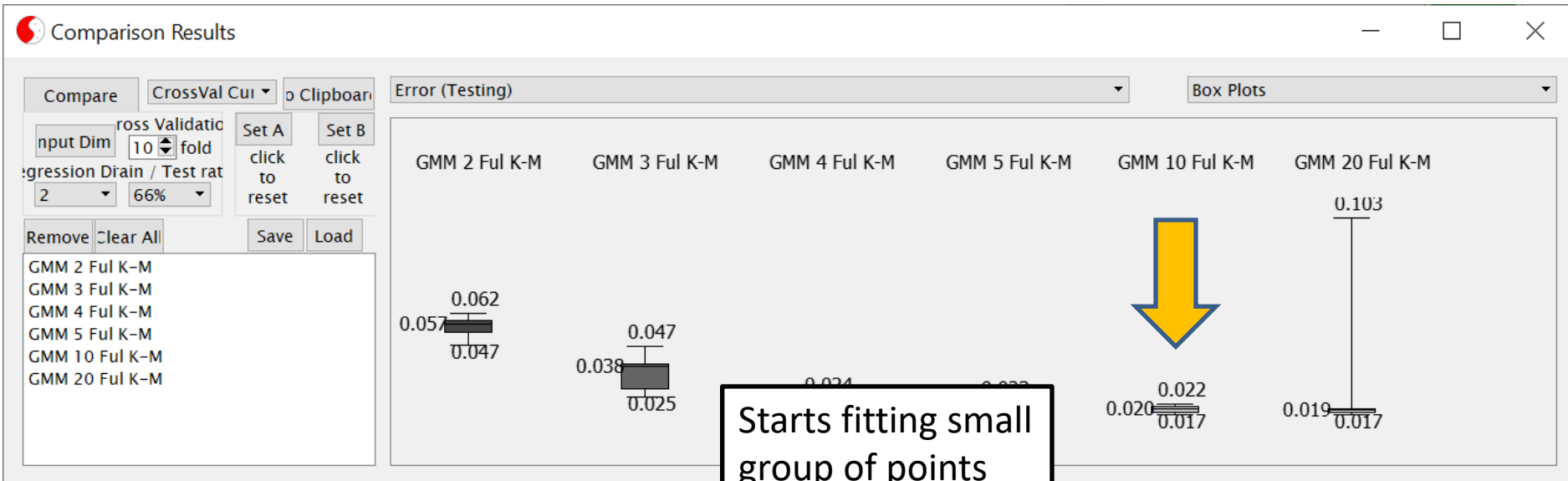
What would be the effect of overfitting in GMR ?

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”.

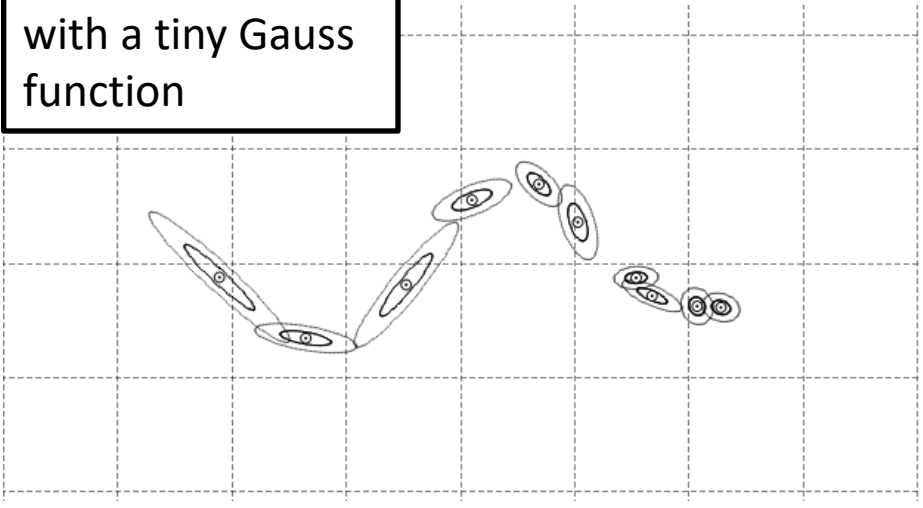
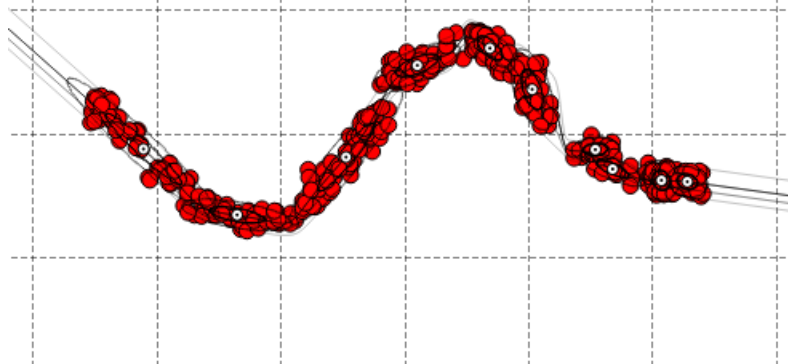


Overfitting with multiple Gauss functions

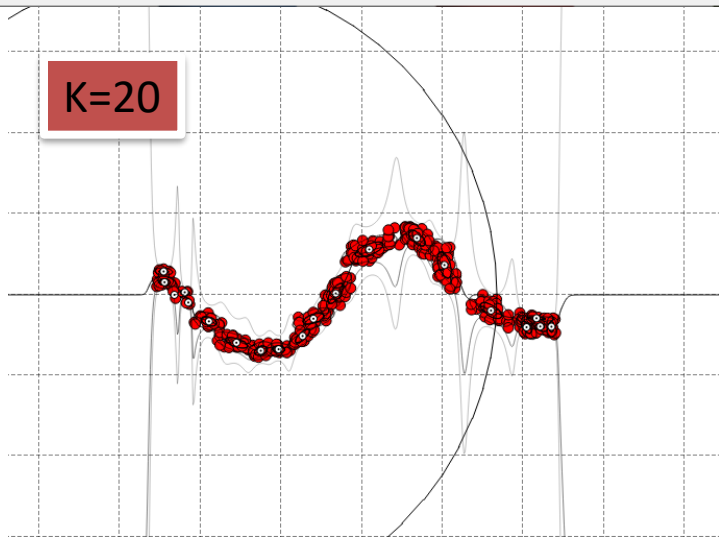
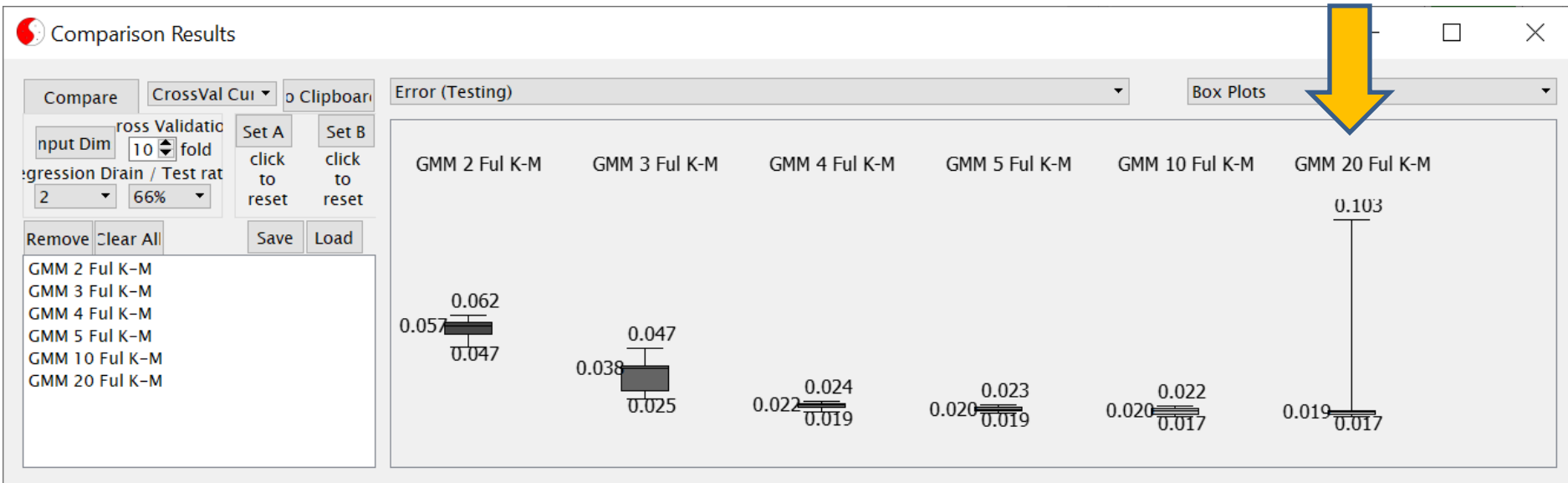


Starts fitting small group of points with a tiny Gauss function

K=10



Overfitting with multiple Gauss functions





Interpolates across two neighbor linear models.

Regressive line is coherent locally, as it is affected only by the influence of neighboring Gauss functions.

Predictions drops away from data



Mean of all local models

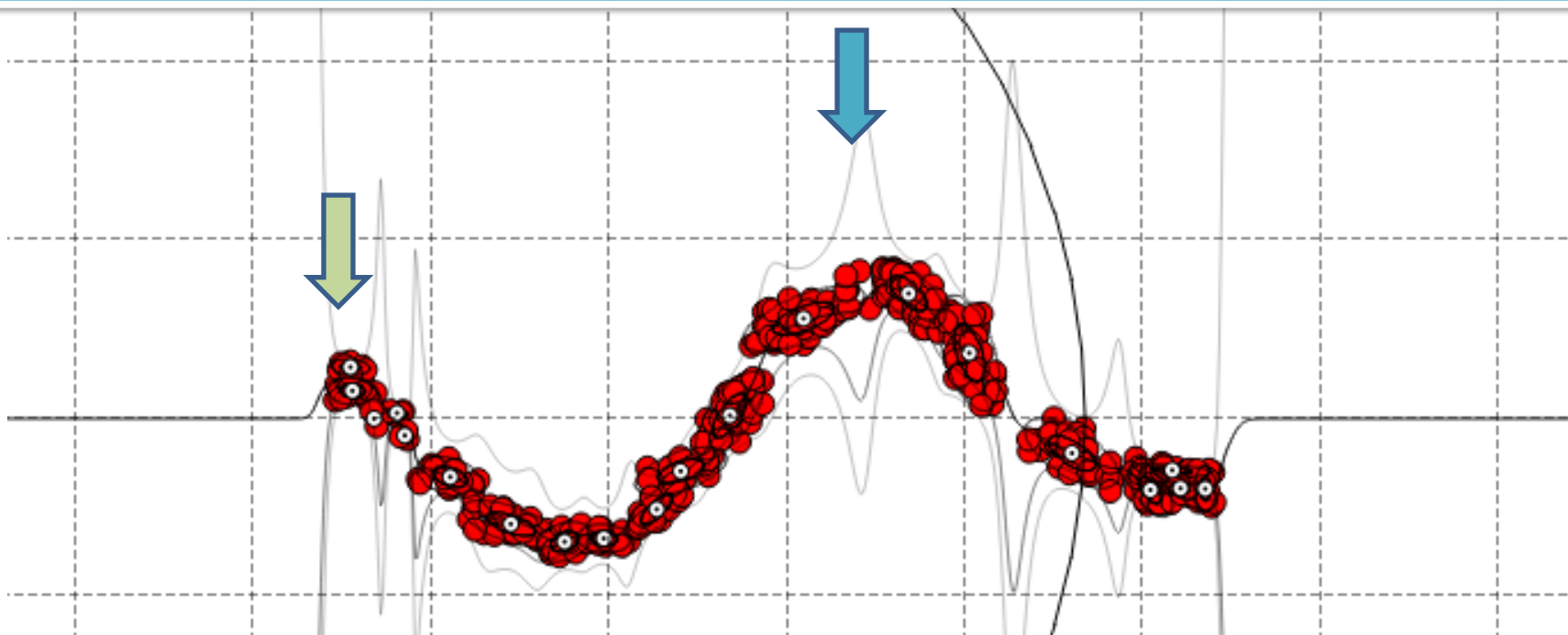
$$y = \sum_{k=1}^K \beta_k(x) \cdot \tilde{\mu}_{y|x}^k(x)$$

~1 for all k

$$\beta_k(x) = \frac{\alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)}{\sum_{k=1}^K \alpha_k \cdot p(x; \mu_x^k, \Sigma_x^k)} \underbrace{\left(\mu_y^k + \Sigma_{yx}^k \left(\Sigma_{xx}^k \right)^{-1} \left(x - \mu_x^k \right) \right)}_{\tilde{\mu}_{y|x}^k(x)}$$

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.

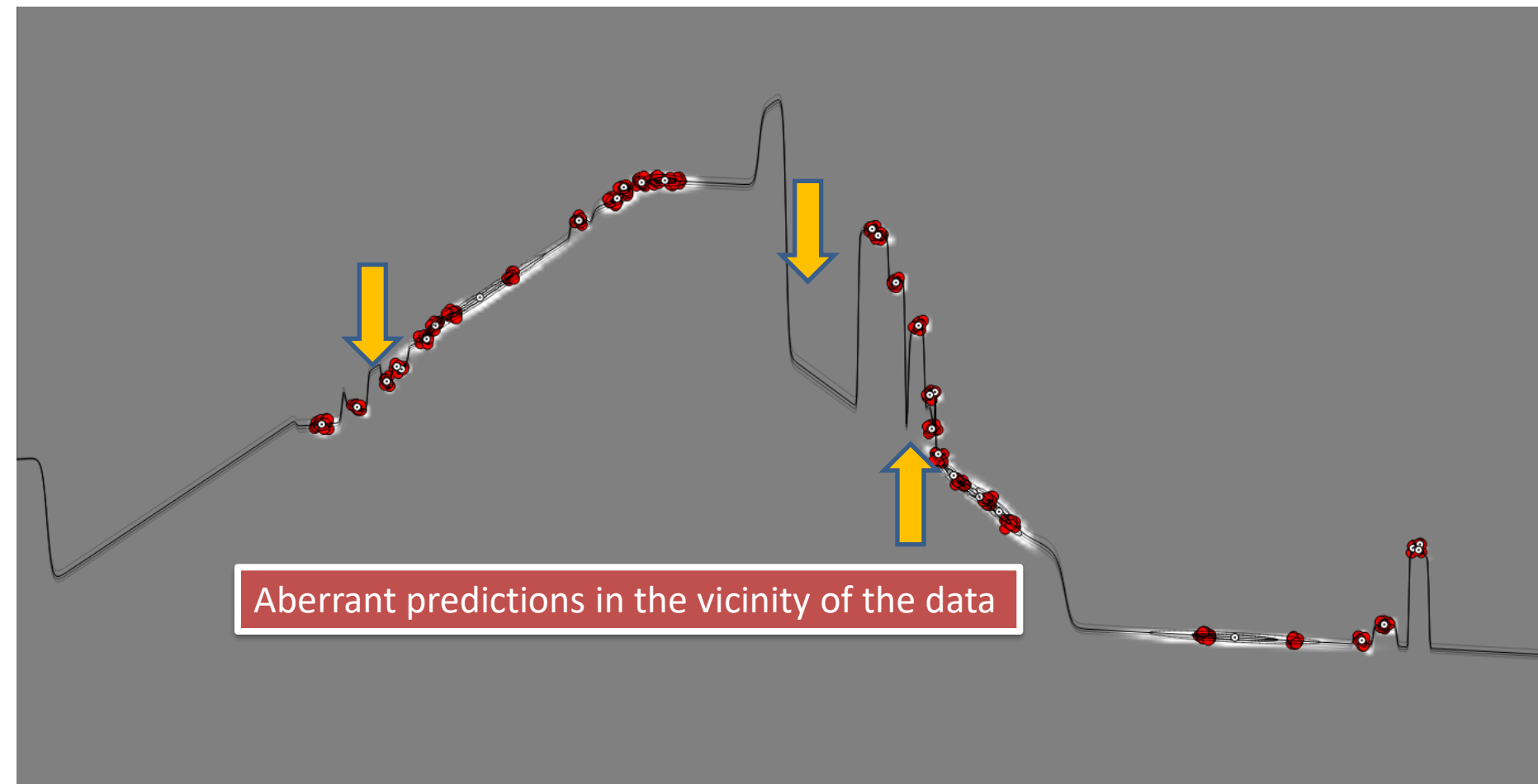


$$\text{var} \{ p(y | x) \} = \sum_{k=1}^K \beta_k(x) \cdot \left(\left(\tilde{\mu}_{y|x}^k(x) \right)^2 + \tilde{\Sigma}_{y|x}^k \right) - \left(\sum_{k=1}^K \left(\beta_k(x) \cdot \tilde{\mu}_{y|x}^k(x) \right) \right)^2$$

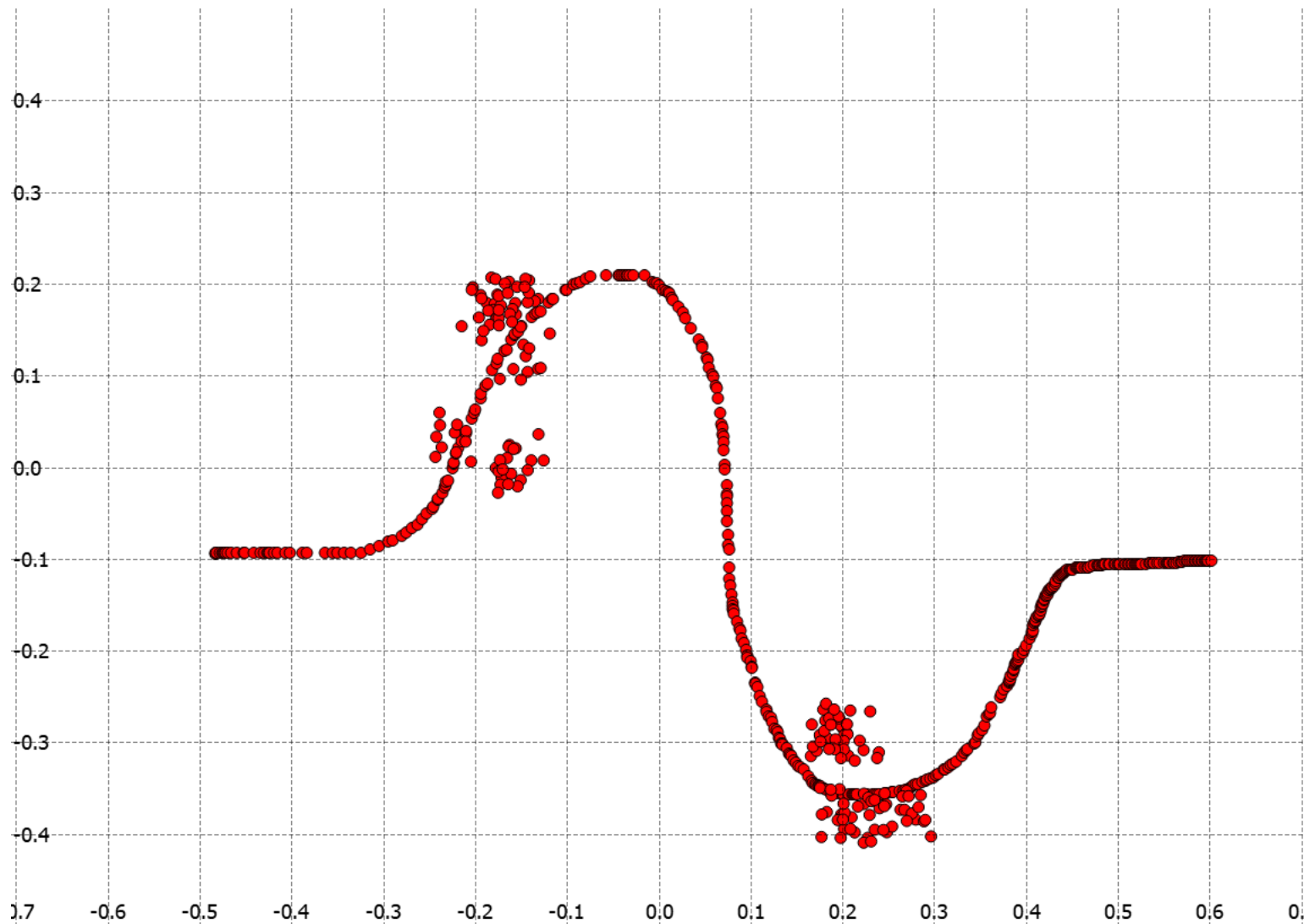
$$\text{with } \tilde{\Sigma}_{y|x}^k = \Sigma_{yy}^k - \Sigma_{yx}^k \left(\Sigma_{xx}^k \right)^{-1} \Sigma_{xy}^k$$

Numerical instabilities

Overfitting with multiple Gauss functions

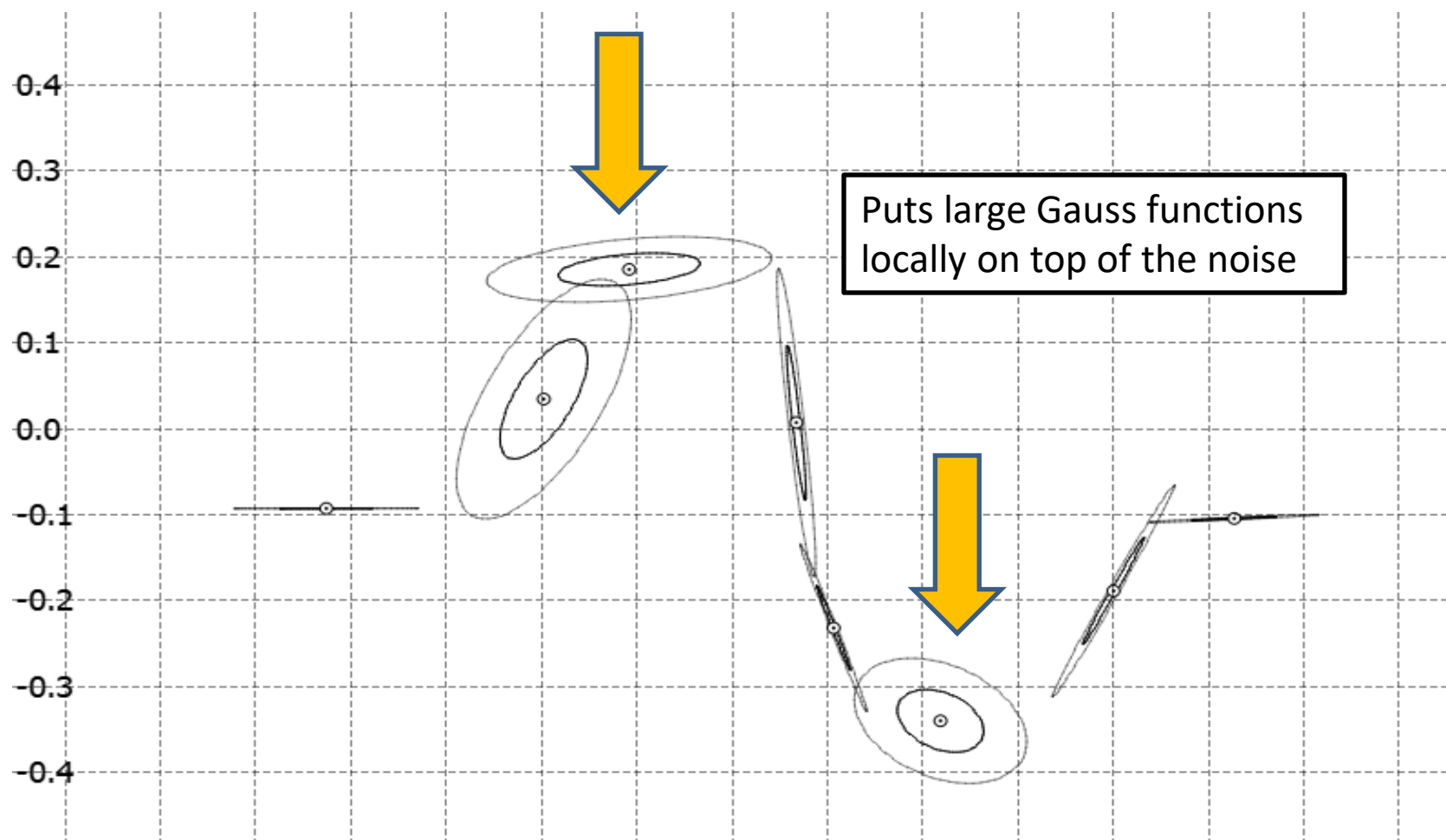


Regression: noise

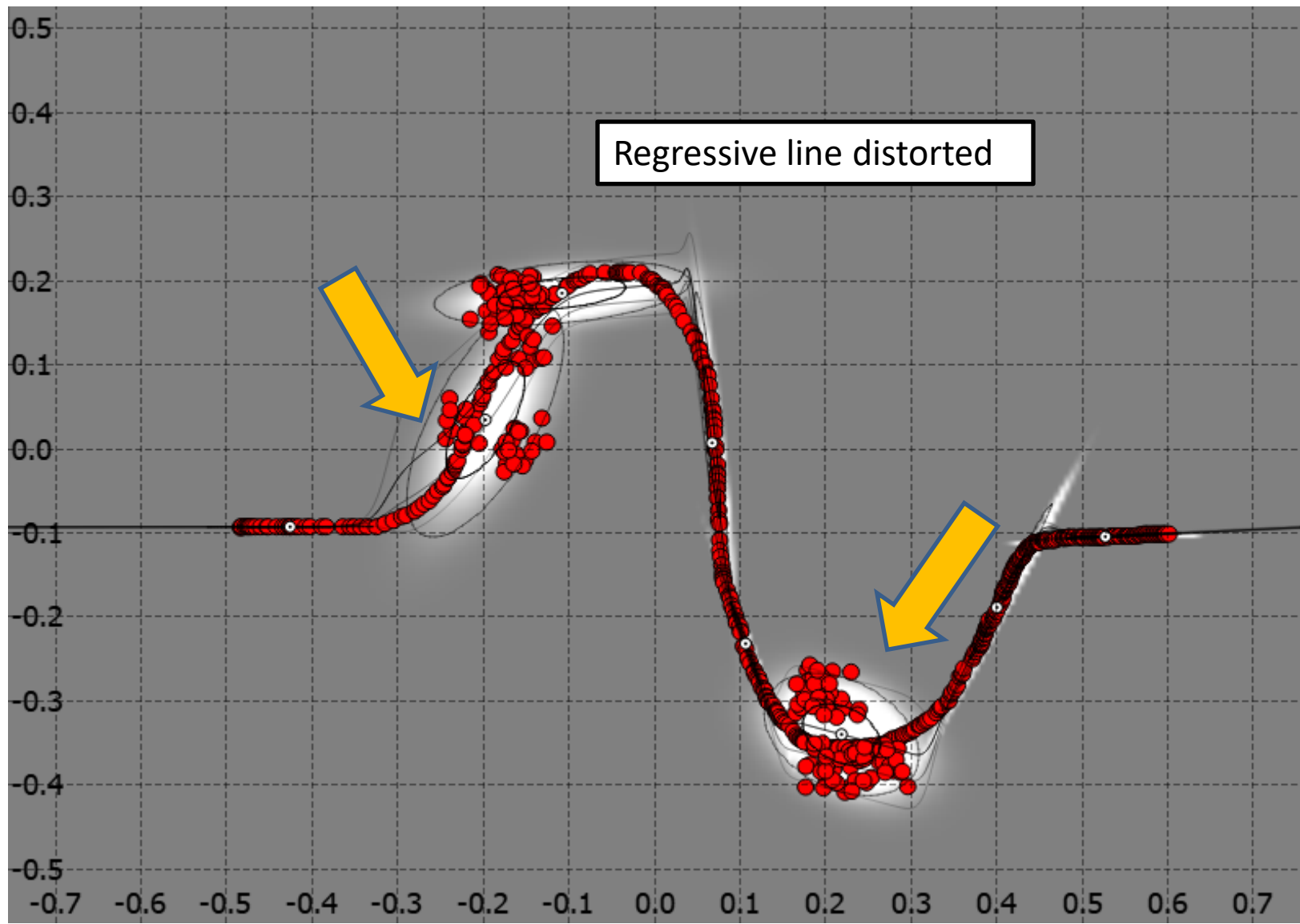


How would GMR handle this noise?

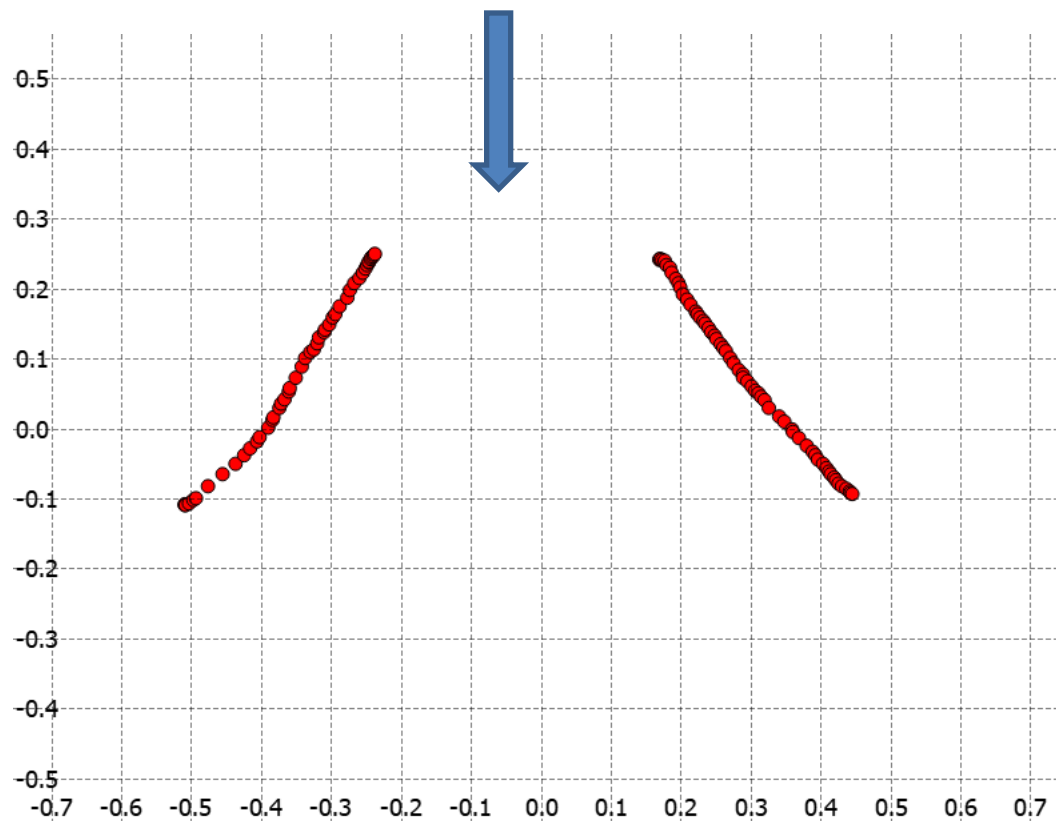
Regression: noise



Regression: noise

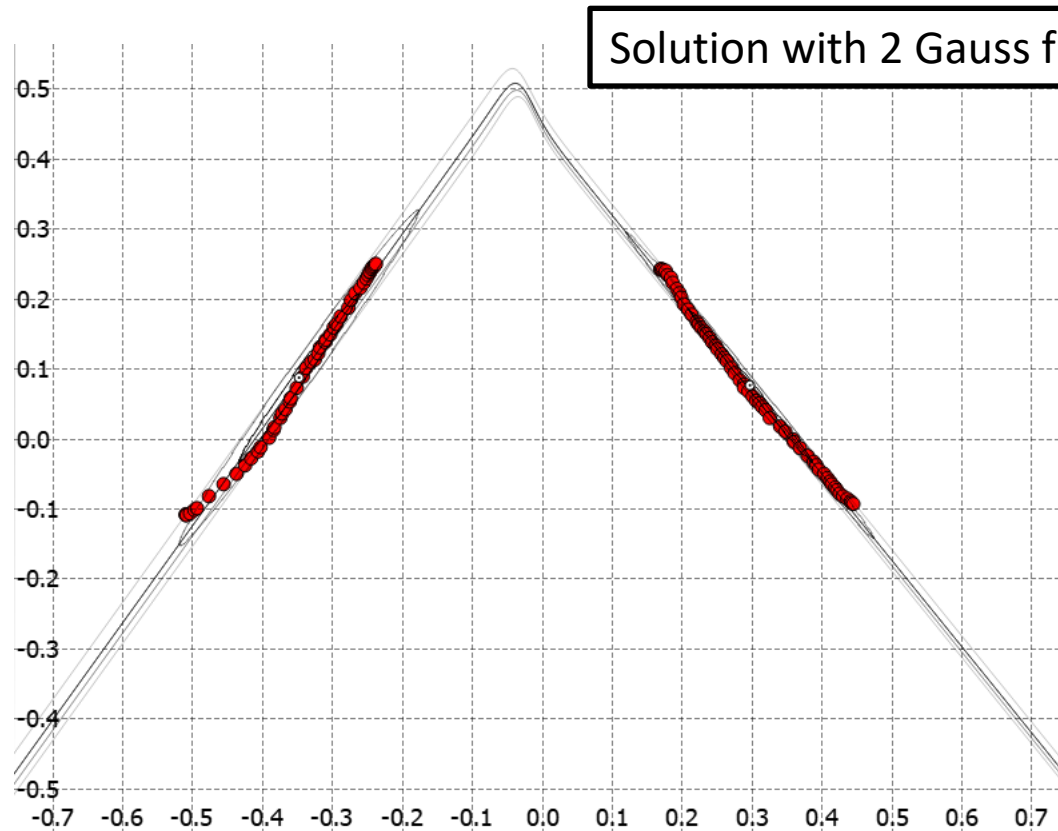


Regression: interpolation



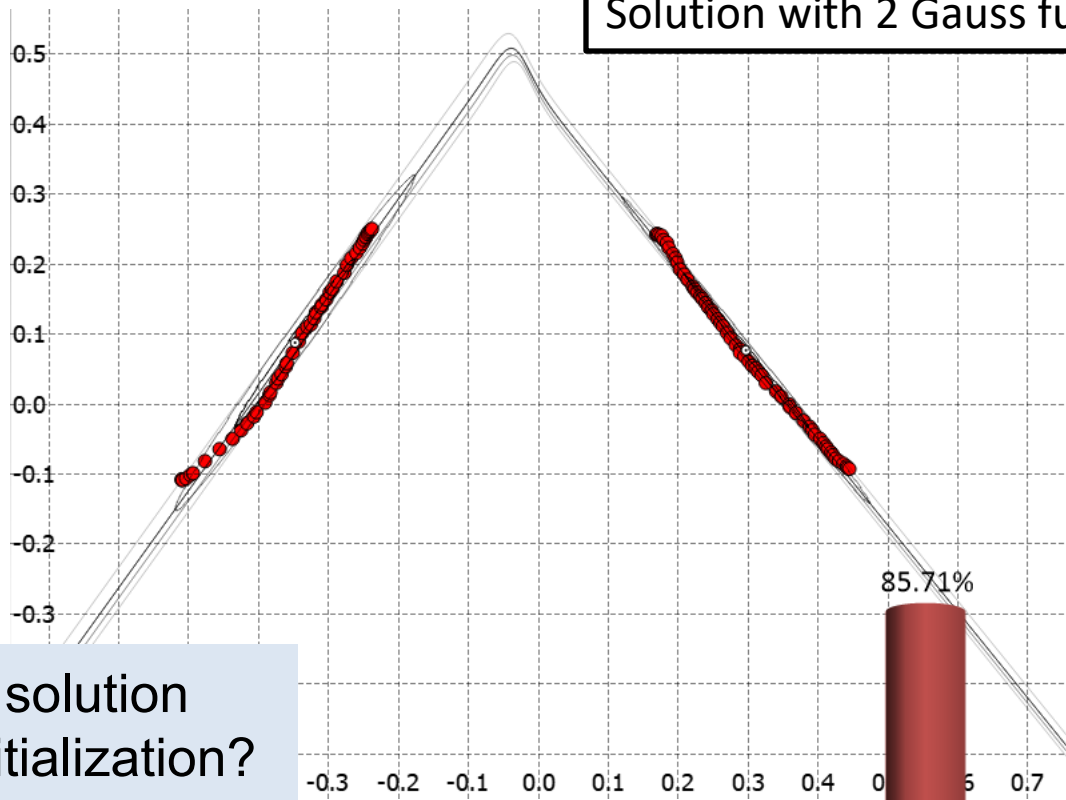
How would GMR handle missing data?

Regression: interpolation



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

Solution with 2 Gauss functions

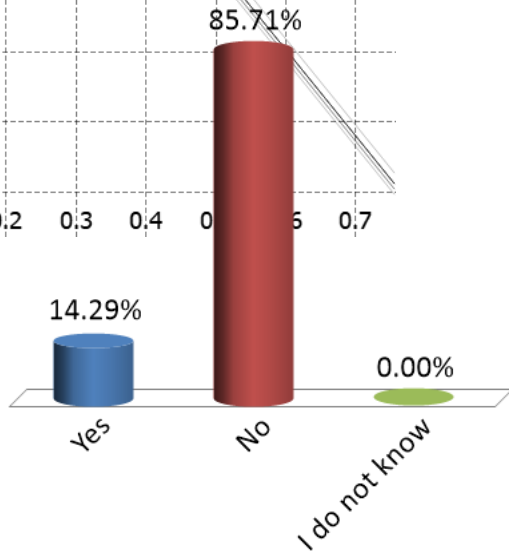


Would the solution depend on initialization?

A. Yes

B. No

C. I do not know

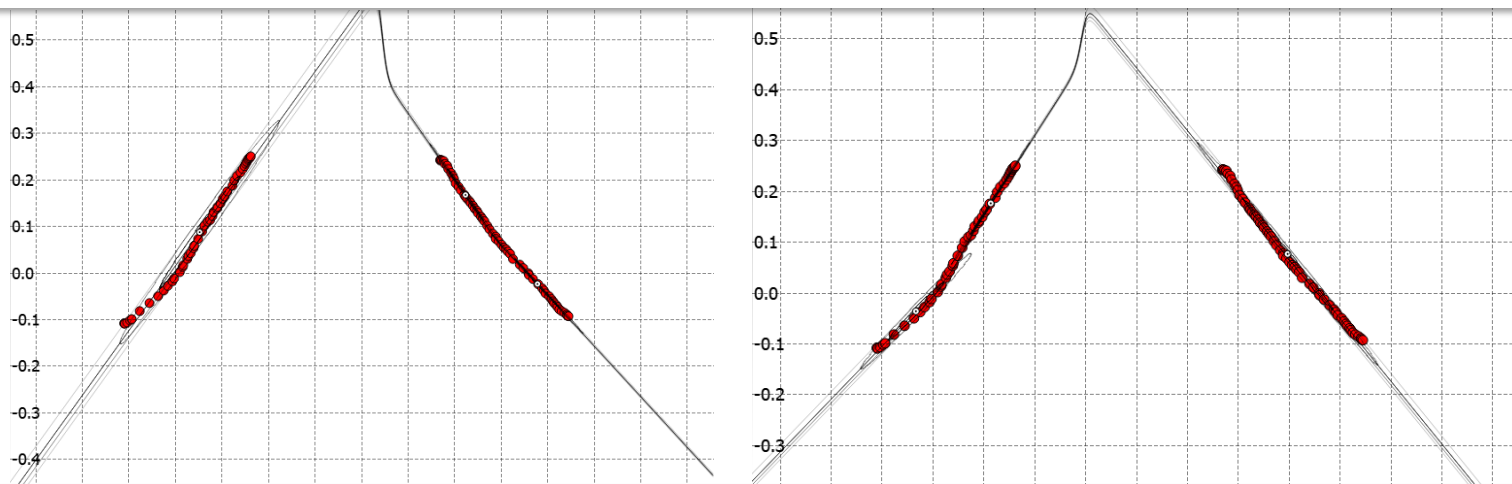


Regression: interpolation



The solution is identical across runs as there is no other local optimum.

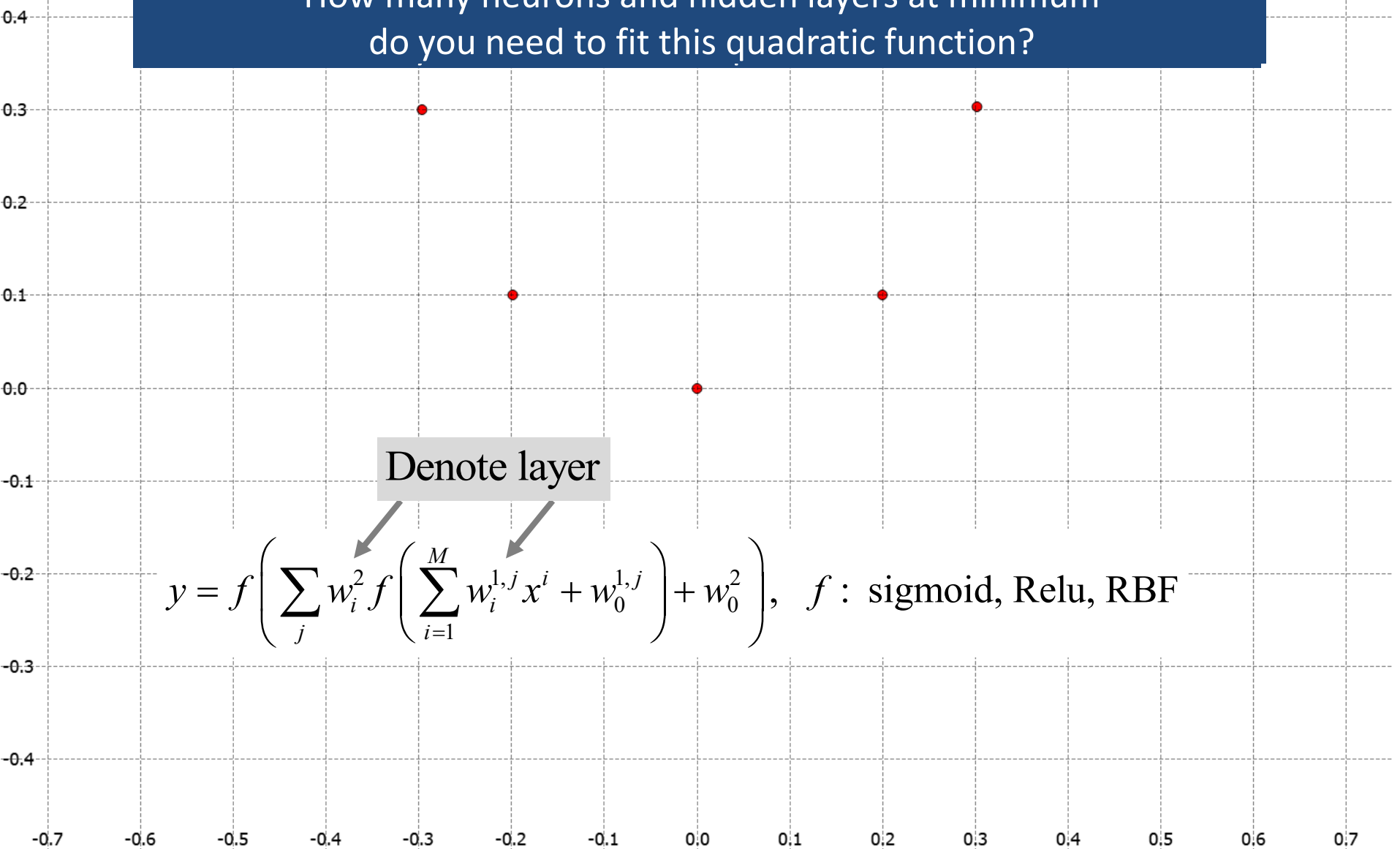
But in principle the solution is not unique and there are often many solutions, each of which corresponds to a local optima on the likelihood.



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

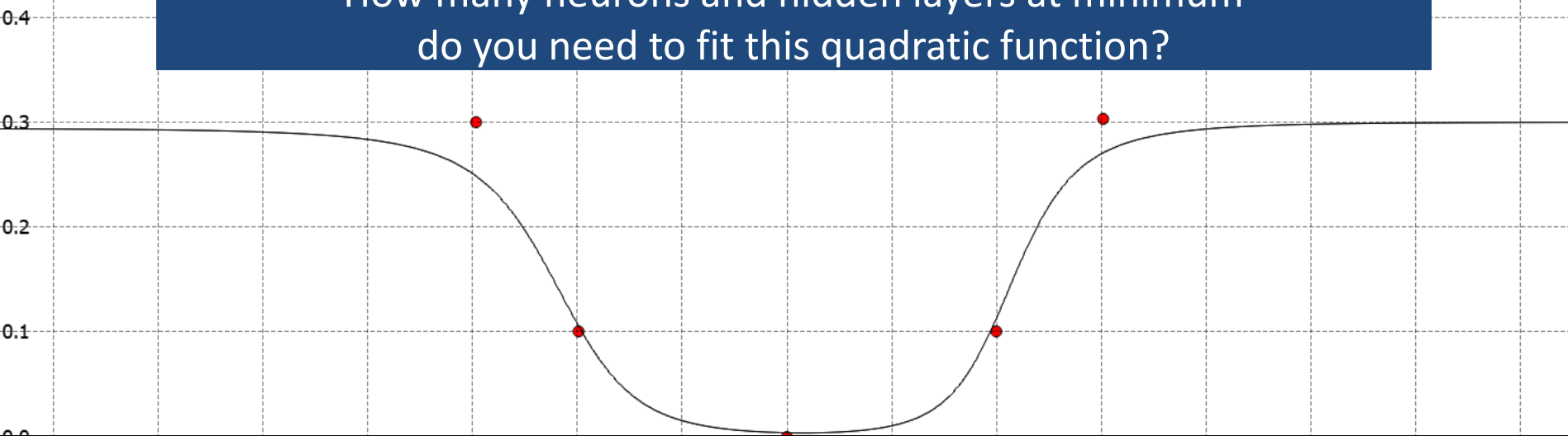
Regression: neural networks

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



Regression: neural networks

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

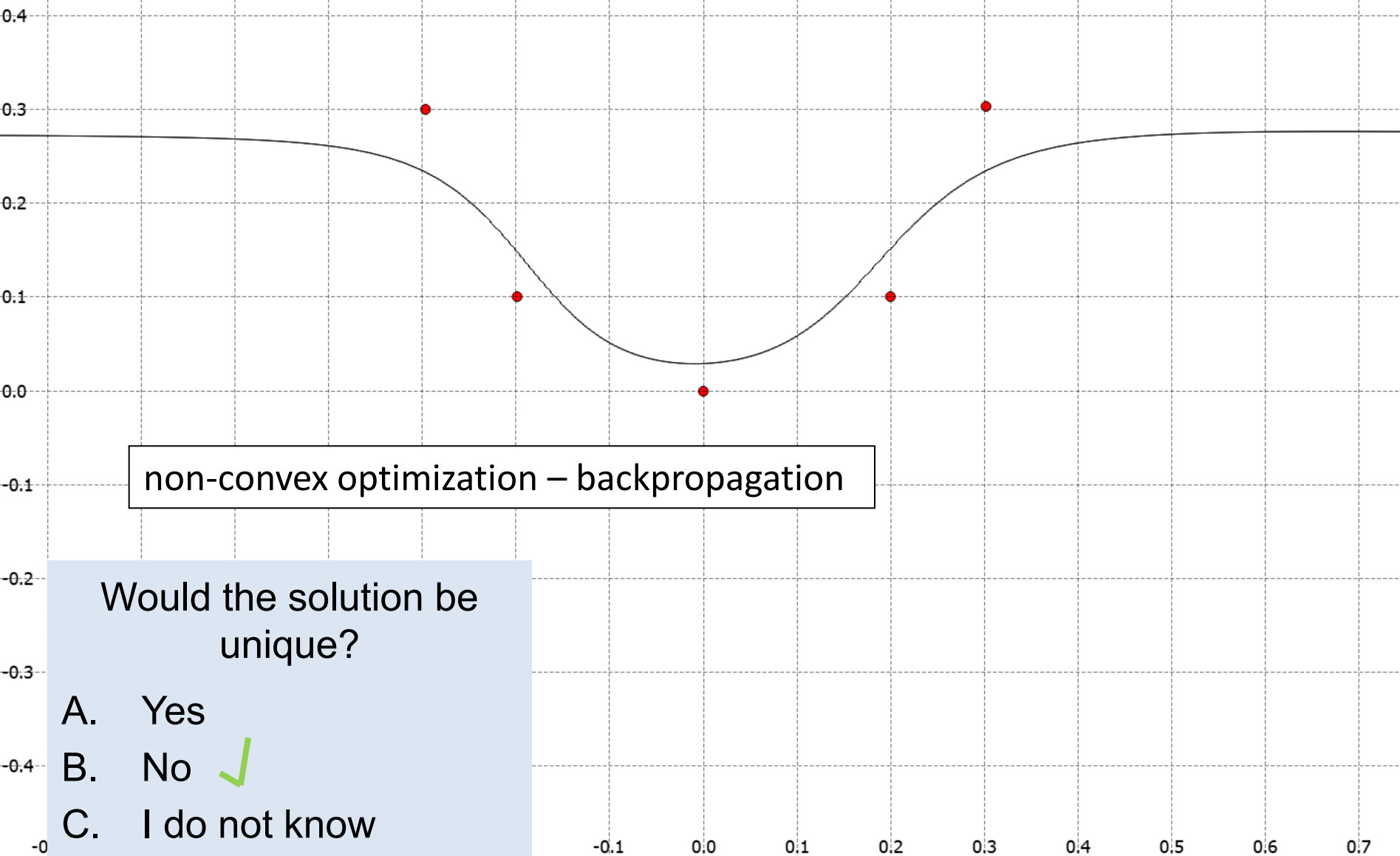


We need 1 hidden layer and at least 2 neurons in the hidden layer (symmetry of the dataset), and 4 to model accurately the change in slope across the 4 points.

$$y = f \left(\sum_j w_j^2 f \left(\sum_{i=1}^M w_i^{1,j} x^i + w_0^{1,j} \right) + w_0^2 \right), \quad f : \text{sigmoid, ReLU}$$

ReLU - **Linear** combination of K local regressive models
 Sigmoid - **Quasi** linear combination

Regression: neural networks



Regression: neural networks

2 neurons in the hidden layer

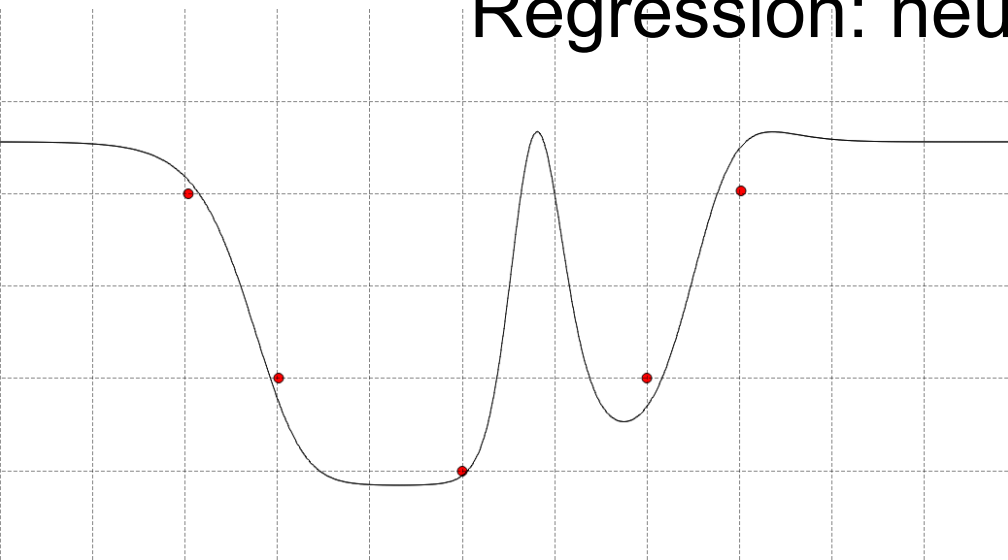
RBF

$$y = f \left(\sum_j w_j^2 \left(f \left(\sum_{i=1}^M w_i^{1,j} x^i + w_0^{1,j} \right) + w_0^2 \right) \right), \quad f : \text{RBF}$$

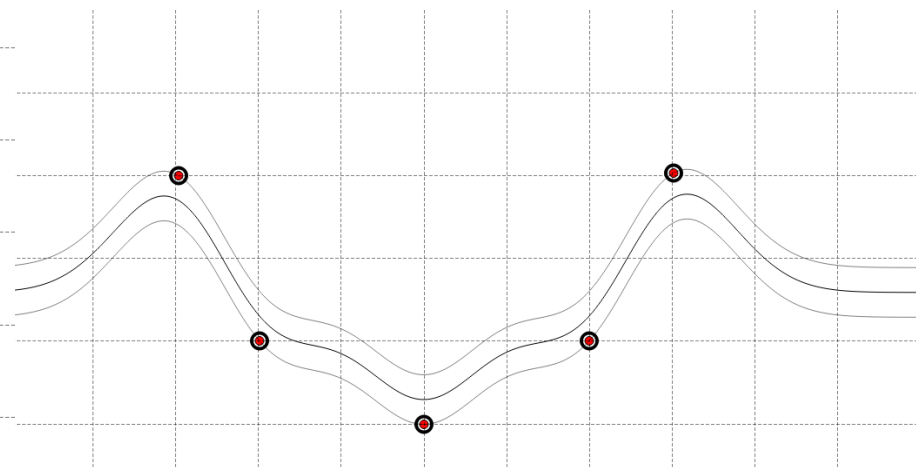
Linear combination of RBF functions – similar expression to SVR

But the RBF are not centered on one point but on a combination of points.

Regression: neural networks



Solution with 4 neurons



SVR solution

$$y = f \left(\sum_j w_j^2 f \left(\sum_{i=1}^M w_i^{1,j} x^i + w_0^{1,j} \right) + w_0^2 \right), \quad f : \text{RBF}$$

RBF function – similar expression to SVR

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