

Support Vector Regression (SVR)



Support Vector Regression: Principle

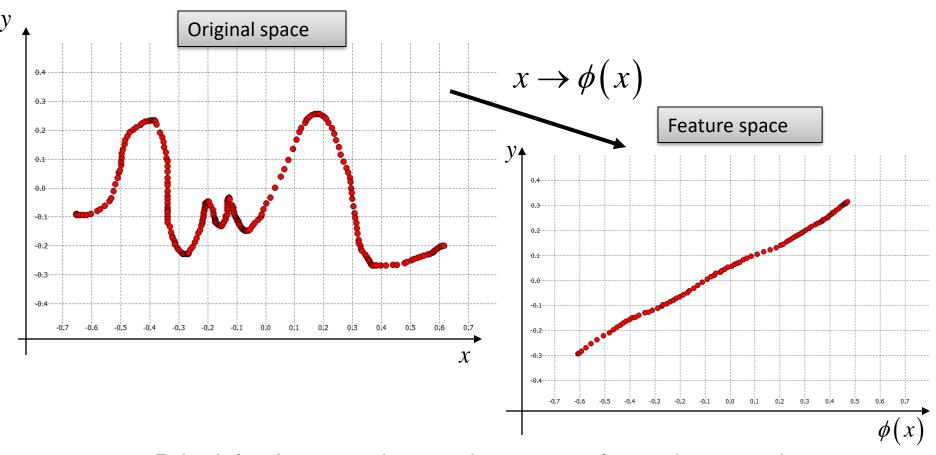
We seek to find a function
$$f$$
, s.t. $y = f(x)$. $y \in \mathbb{R}, x \in \mathbb{R}^N$

Generalize the support vector machine framework for classification to estimate continuous functions

Assume a non-linear mapping through feature space and then perform *linear regression* in feature space.



Support Vector Regression: Principle



→ Perform linear regression in feature space

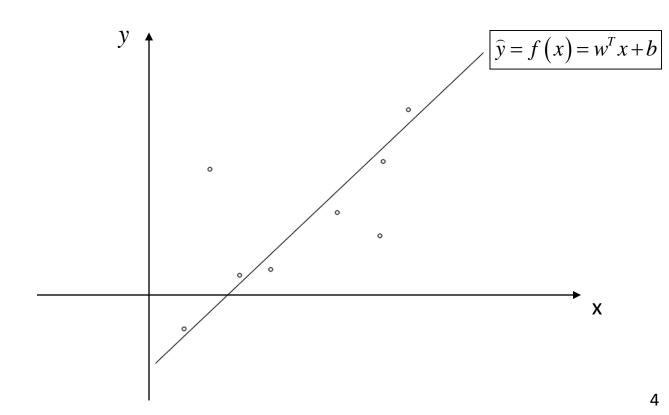


Support Vector Regression: Linear Case

Assume a linear mapping f, s.t. $y = f(x) = w^{T}x + b$.

 $x \in \mathbb{R}^N, y \in \mathbb{R}$

Need to estimate w and b to best predict the pair of training points $\{x^i, y^i\}_{i=1,\dots,M}$.

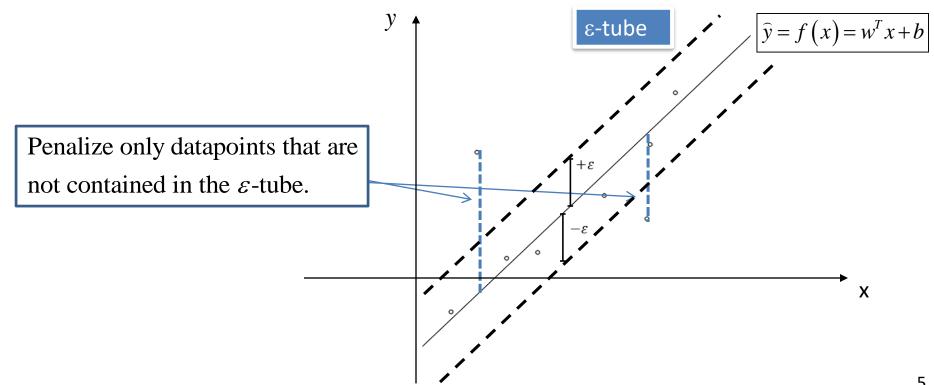




Support Vector Regression: ε-tube

Assume deterministic noise model: $y \pm \varepsilon$

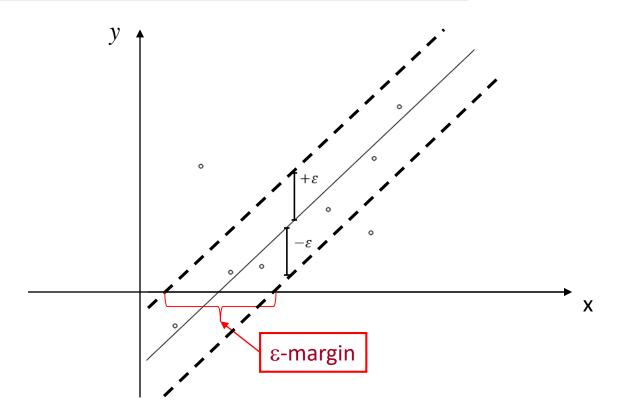
Consider as correctly fit all points such that $f(x) - y \le \varepsilon$.





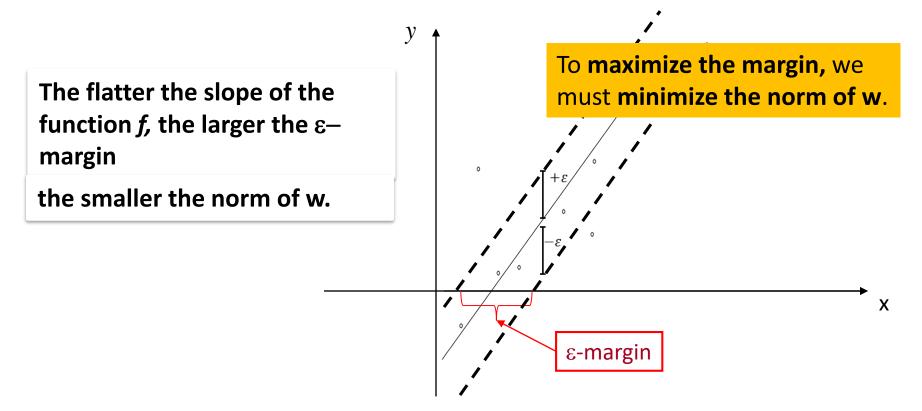
How to assess how well we do for a choice of ε -tube.

The ε -margin is a measure of the width of the ε -insensitive tube. It is a measure of the precision of the regression.





The ε -margin is a measure of the width of the ε -insensitive tube. It is a measure of the precision of the regression.





This can be rephrased as a constraint-based optimization problem:

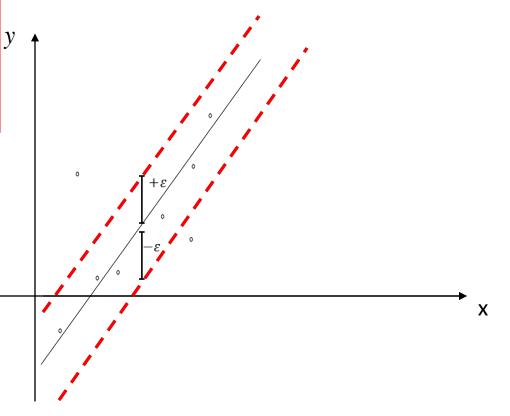
minimize
$$\frac{1}{2} \|w\|^2$$

subject to
$$\begin{cases} \langle w, x^i \rangle + b - y^i \le \varepsilon \\ y^i - \langle w, x^i \rangle - b \le \varepsilon \end{cases}$$

 $\forall i = 1,...M$

Consider as correctly fit all

points such that $|f(x) - y| \le \varepsilon$.





This can be rephrased as a constraint-based optimization problem:

minimize
$$\frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^{M} (\xi_i + \xi_i^*), \quad C > 0$$

subject to
$$\begin{cases} \left\langle w, x^{i} \right\rangle + b - y^{i} \leq \varepsilon + \xi_{i}^{*} \\ y^{i} - \left\langle w, x^{i} \right\rangle - b \leq \varepsilon + \xi_{i} \end{cases}$$

$$\forall i = 1, ... M \qquad \xi_{i} \geq 0, \quad \xi_{i}^{*} \geq 0$$

Introduce slack variables ξ_i, ξ_i^*

Need to penalize points outside the ϵ -insensitive tube.



Support Vector Regression: optimization

We can solve this quadratic problem by introducing sets of α , $\eta \in \mathbb{R}$ Lagrange multipliers and writing the Lagrangian :

Lagrangian = Objective function + multipliers * constraints

$$L(w, \xi, \xi^*, b) = \frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^{M} (\xi_i + \xi_i^*) - \frac{C}{M} \sum_{i=1}^{M} (\eta_i \xi_i + \eta_i^* \xi_i^*)$$
$$-\sum_{i=1}^{M} \alpha_i (\varepsilon + \xi_i - y^i + \langle w, x^i \rangle + b)$$
$$-\sum_{i=1}^{M} \alpha_i^* (\varepsilon + \xi_i^* + y^i - \langle w, x^i \rangle - b)$$



Support Vector Regression: solution

Requiring that the partial derivatives are all zero:

$$\frac{\partial \mathbf{L}}{\partial b} = \sum_{i=1}^{M} \left(\alpha_i - \alpha_i^* \right) = 0. \qquad \longrightarrow \sum_{i=1}^{M} \alpha_i = \sum_{i=1}^{M} \alpha_i^*$$

$$ightarrow \sum_{i=1}^{M} lpha_i = \sum_{i=1}^{M} lpha_i^*$$

Rebalancing the effect of the support vectors on both sides of the ε -tube

$$\frac{\partial \mathbf{L}}{\partial w} = w - \sum_{i=1}^{M} (\alpha_i - \alpha_i^*) x^i = 0.$$

$$\Rightarrow w = \sum_{i=1}^{M} (\alpha_i - \alpha_i^*) x^i.$$

Linear combination of support vectors

$$y = f(x)$$

$$= \langle w, x \rangle + b$$

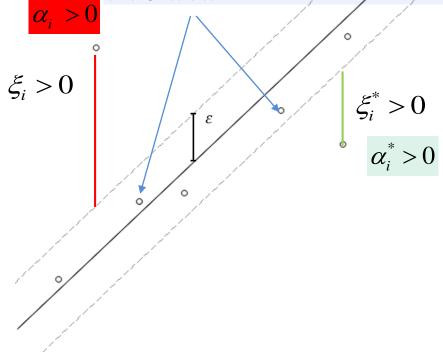
$$= \sum_{i=1}^{M} (\alpha_i - \alpha_i^*) \langle x^i, x \rangle + b$$

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



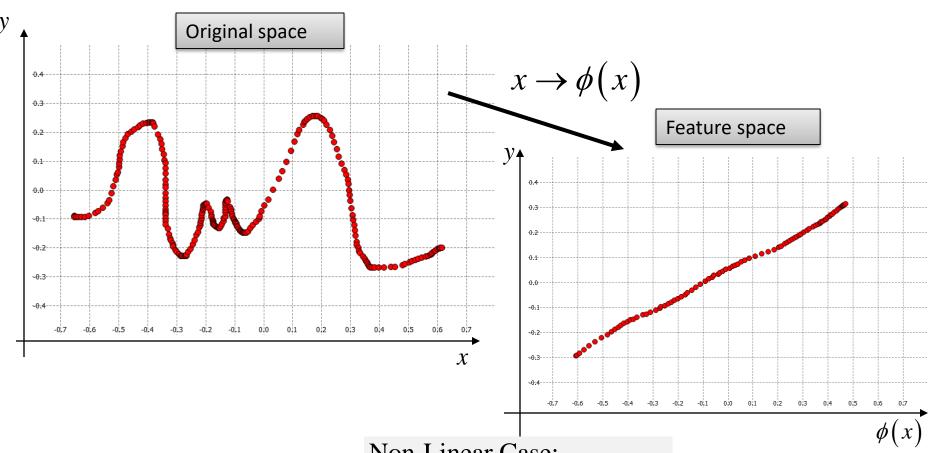
Support Vector Regression: solution





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





Linear Case:

$$y = f(x) = \langle w, x \rangle + b$$

Non-Linear Case:

$$x \to \phi(x)$$

w lives in feature space!

$$y = f(\phi(x)) = \langle w, \phi(x) \rangle + b$$



Replacing in the primal Lagrangian, we get the Dual optimization:

$$\max_{\alpha,\alpha^*} \begin{cases} -\frac{1}{2} \sum_{i,j=1}^{M} (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) \left\langle \phi(x^i), \phi(x^j) \right\rangle \\ -\varepsilon \sum_{i=1}^{M} (\alpha_i^* + \alpha_i) + \sum_{i=1}^{M} y^i (\alpha_i^* + \alpha_i) \end{cases}$$
subject to
$$\sum_{i=1}^{M} (\alpha_i^* - \alpha_i) = 0 \text{ and } \alpha_i^*, \alpha_i^i \in [0, C]$$

Kernel Trick

$$k(x^{i}, x^{j}) = \langle \phi(x^{i}), \phi(x^{j}) \rangle$$



The solution is given by:

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

An estimate of *b* can be computed from the KKT conditions:

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

using only the SVs on the ε -tube.



The solution is given by:

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

Linear Coefficients (Lagrange multipliers for each constraint).

If one uses RBF Kernel,

M un-normalized isotropic

Gaussians centered on each training datapoint.

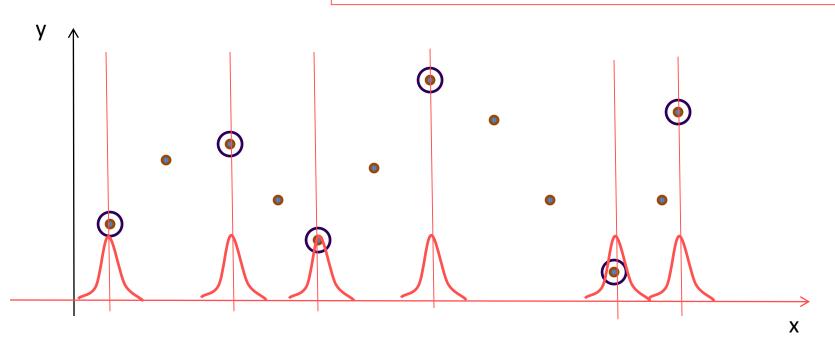


Support Vector Regression: interpretation

The solution is given by:

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

Kernel places a Gauss function on each SV





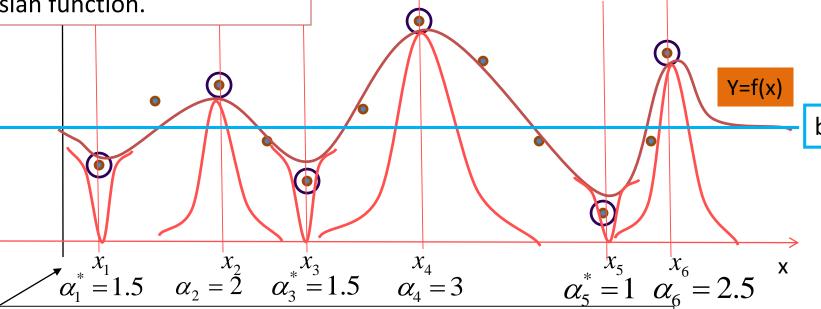
Support Vector Regression: interpretation

The solution is given by:

Converges to b when SV effect vanishes.

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

The Lagrange multipliers define the importance of each Gaussian function.





ε-SVR: 3 Hyperparameters

The solution to SVR we just saw is referred to as ϵ –SVR

Two Hyperparameters for optimization

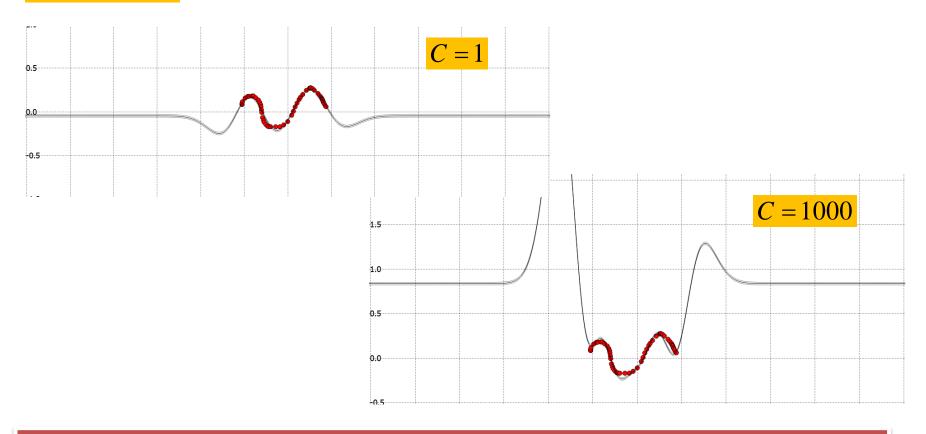
minimize
$$\frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^{M} (\xi_i + \xi_i^*)$$

$$\begin{cases} \langle w, x^i \rangle + b - y^i \le \varepsilon + \xi_i^* \\ y^i - \langle w, x^i \rangle - b \le \varepsilon + \xi_i \\ \xi_i \ge 0, \quad \xi_i^* \ge 0 \end{cases}$$
subject to

- 1. C controls the penalty term on poor fit.
- 2. ε determines the minimal required precision for the fit.
- 3. The kernel width for the RBF kernel determines the locality of the regression

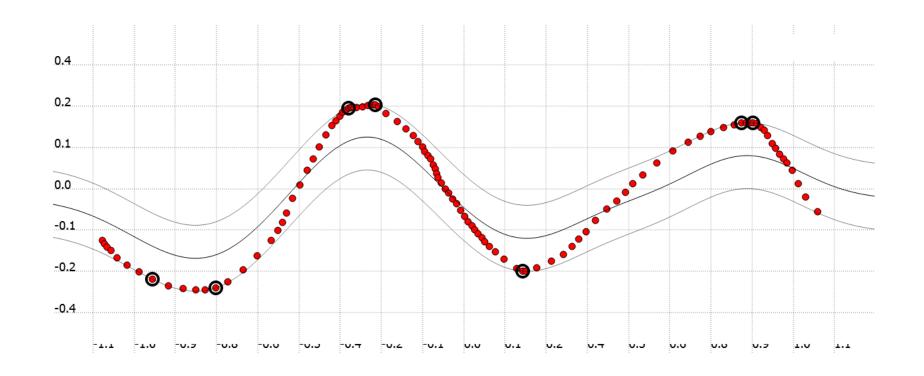


$$\alpha_i, \alpha_i^* \in [0, C]$$



Effect of C on the fit.

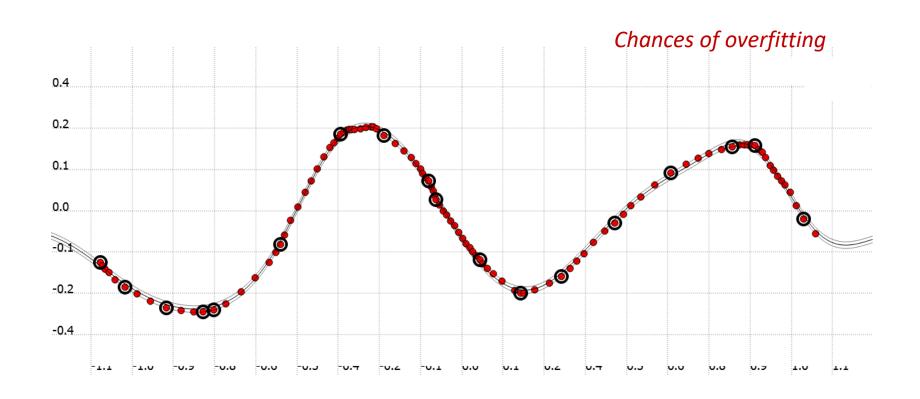




Here fit using C=1, ε =0.1, kernel width=0.01.

Effect of the ε on the fit.

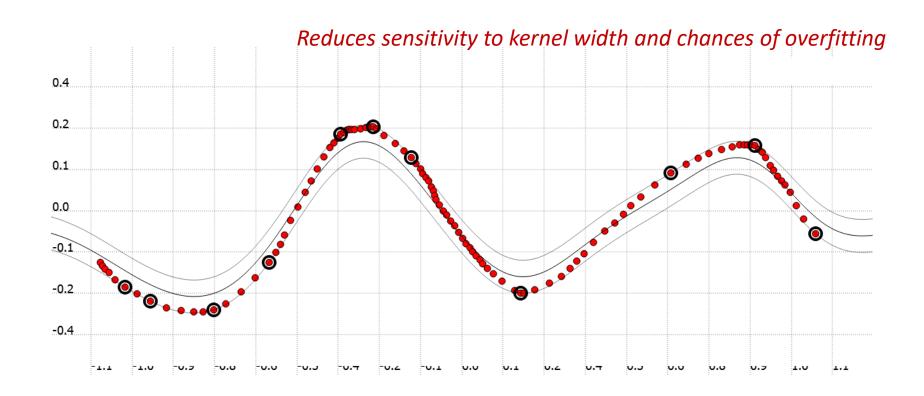




Here fit using C=1, ε =0.01, kernel width=0.01

Effect of the ε on the fit.

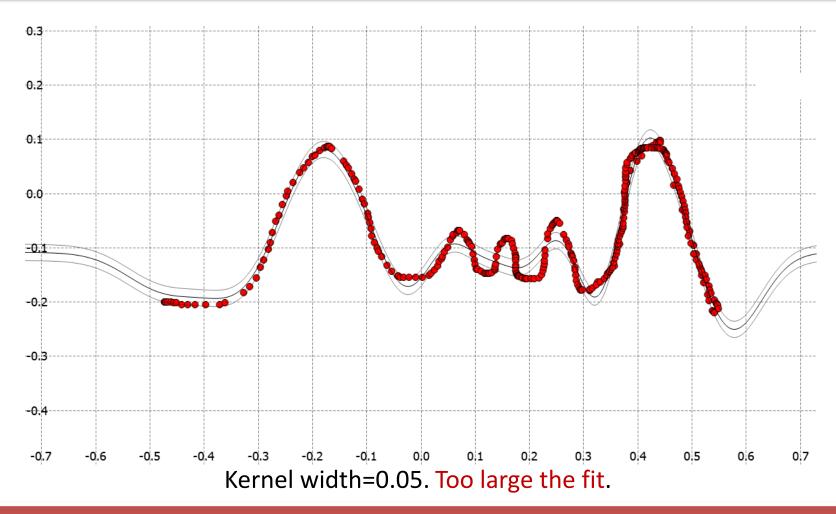




Here fit using C=1, ε =0.05, kernel width=0.01

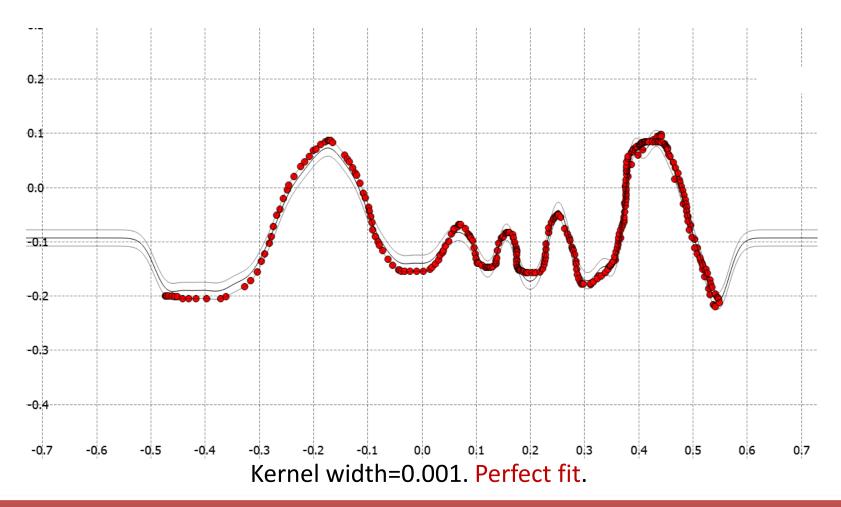
Effect of the ε on the fit.





Effect of the RBF kernel width on the fit.





Effect of the RBF kernel width on the fit.



Summary

<u>Linear regression</u> can be solved through Least-Mean-Square estimation and yields an optimal analytical solution.

Weighted regression offers the possibility to perform a local regression and yields also an optimal analytical solution.

The estimate is no longer global and is computed around each group of data point!

<u>Support Vector Regression</u>: performs regression on a non-linear function. Determines automatically the important points. The estimate is globally optimal.