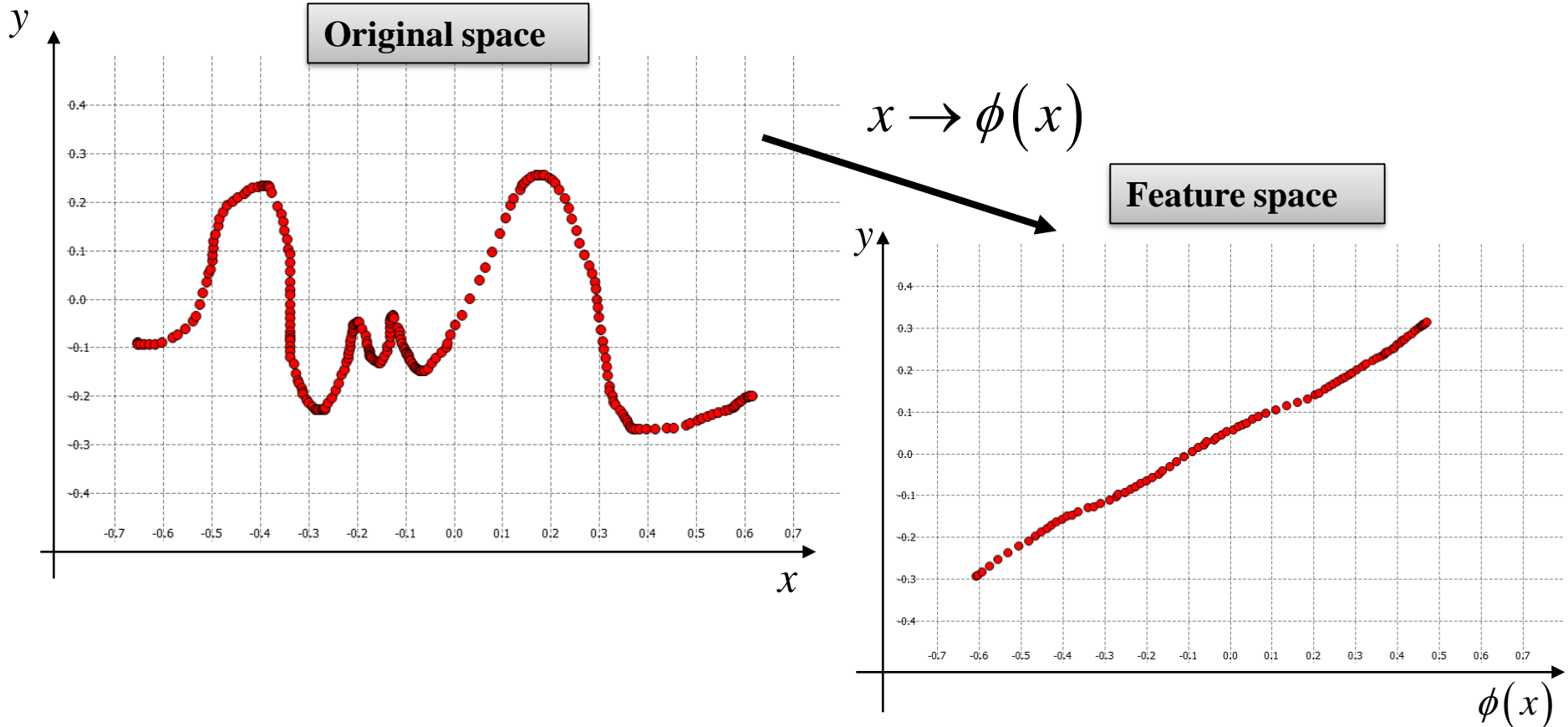


***Kernel trick***

***SVR extensions***  
 ***$\nu$ -SVR, RVR***



# Recap: Nonlinear regression with SVR



Principle: Assume there exists a transformation  $\phi$  such that the problem becomes linear

→ Perform linear regression in feature space



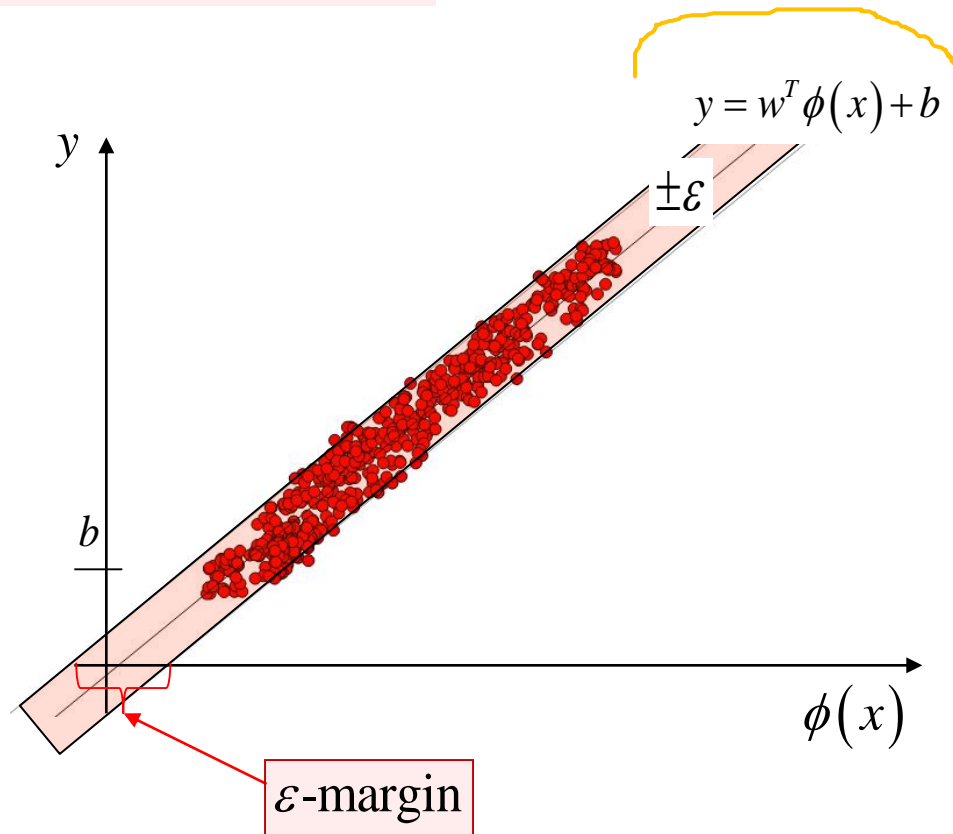
# Recap: Nonlinear regression with SVR

The  $\varepsilon$ -insensitive tube is the region surrounding the regression function that encapsulates acceptable noise.

To maximize the  $\varepsilon$ -margin, minimize the norm of  $w$ .

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

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



# Recap: Nonlinear regression with SVR

This can be expressed as a constraint-based optimization problem of the form:

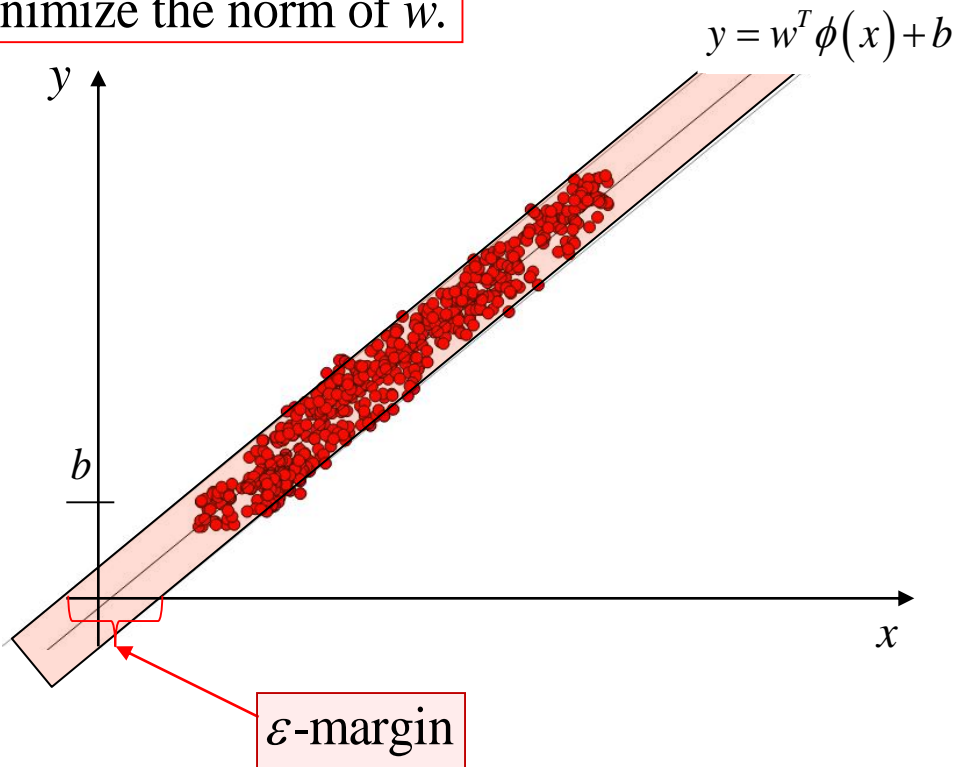
minimize  $\frac{1}{2} \|w\|^2$

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

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

maximize the  $\varepsilon$ -margin,  
minimize the norm of  $w$ .

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

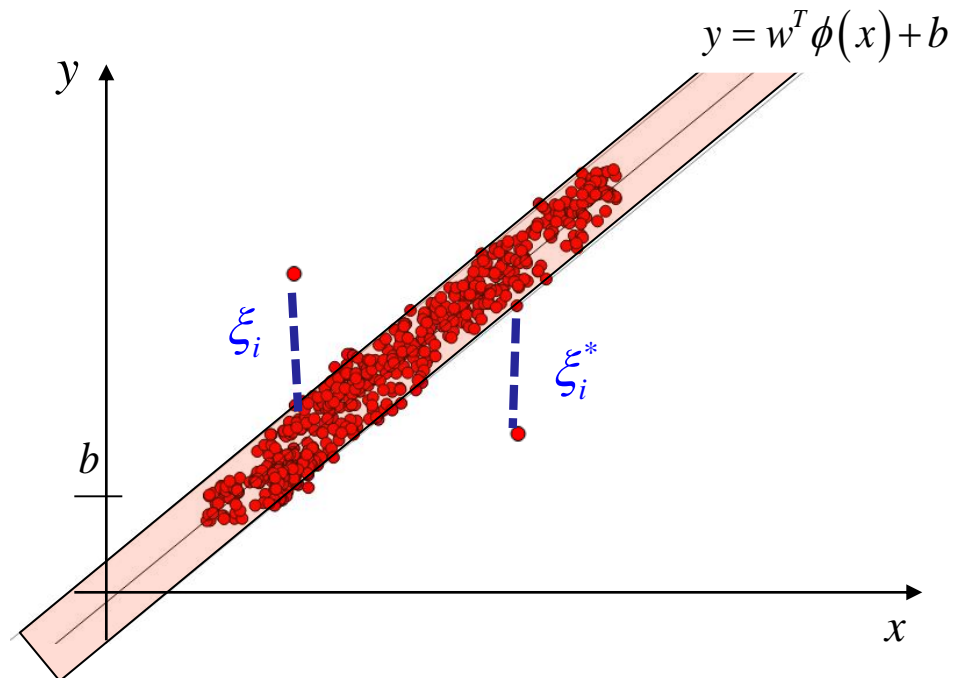


# Recap: Nonlinear regression with SVR

Relaxes constraints and introduce slack variables  $\xi_i, \xi_i^*$ .

Penalizes for points outside the  $\varepsilon$ -tube,  $C \geq 0$

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^M (\xi_i + \xi_i^*) \\ &\text{subject to} \quad \begin{cases} \langle w, \phi(x^i) \rangle + b - y^i \leq \varepsilon + \xi_i \\ y^i - \langle w, \phi(x^i) \rangle - b \leq \varepsilon + \xi_i^* \\ \xi_i \geq 0, \quad \xi_i^* \geq 0 \end{cases} \end{aligned}$$



# Recap: SVR Hyperparameters

$$\begin{aligned}
 &\text{minimize } \frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^M (\xi_i + \xi_i^*) \\
 &\text{subject to } \begin{cases} \langle w, \phi(x^i) \rangle + b - y^i \leq \varepsilon + \xi_i \\ y^i - \langle w, \phi(x^i) \rangle - b \leq \varepsilon + \xi_i^* \\ \xi_i \geq 0, \quad \xi_i^* \geq 0 \end{cases}
 \end{aligned}$$

$$\begin{aligned}
 y &= f(x) \\
 &= \sum_{i=1}^M \alpha_i k(x^i, x) + b
 \end{aligned}$$

**Dual optimization problem:**

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



# Recap: SVR Hyperparameters

Value for hyperparameter C is unbounded.  $C > 0$ .

$$\begin{aligned} &\text{minimize } \frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^M (\xi_i + \xi_i^*) \\ &\text{subject to } \begin{cases} \langle w, \phi(x^i) \rangle + b - y^i \leq \varepsilon + \xi_i \\ y^i - \langle w, \phi(x^i) \rangle - b \leq \varepsilon + \xi_i^* \\ \xi_i \geq 0, \quad \xi_i^* \geq 0 \end{cases} \end{aligned}$$

**Dual optimization problem:**

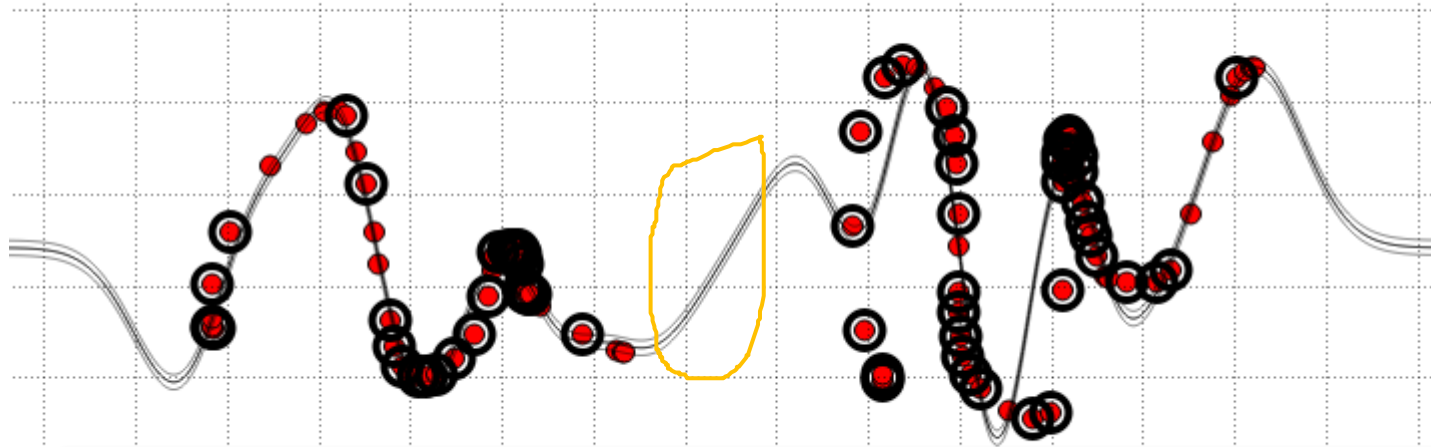
$$\max_{\alpha, \alpha^*} \begin{cases} -\frac{1}{2} \sum_{i,j=1}^M \alpha_i \alpha_j \cdot \langle \phi(x^i), \phi(x^j) \rangle \\ -\varepsilon \sum_{i=1}^M \alpha_i + \sum_{i=1}^M y^i \alpha_i \end{cases}$$

$$\text{subject to } \sum_{i=1}^M \alpha_i = 0 \quad \text{and} \quad \alpha_i \in \left[ 0, \frac{C}{M} \right]$$

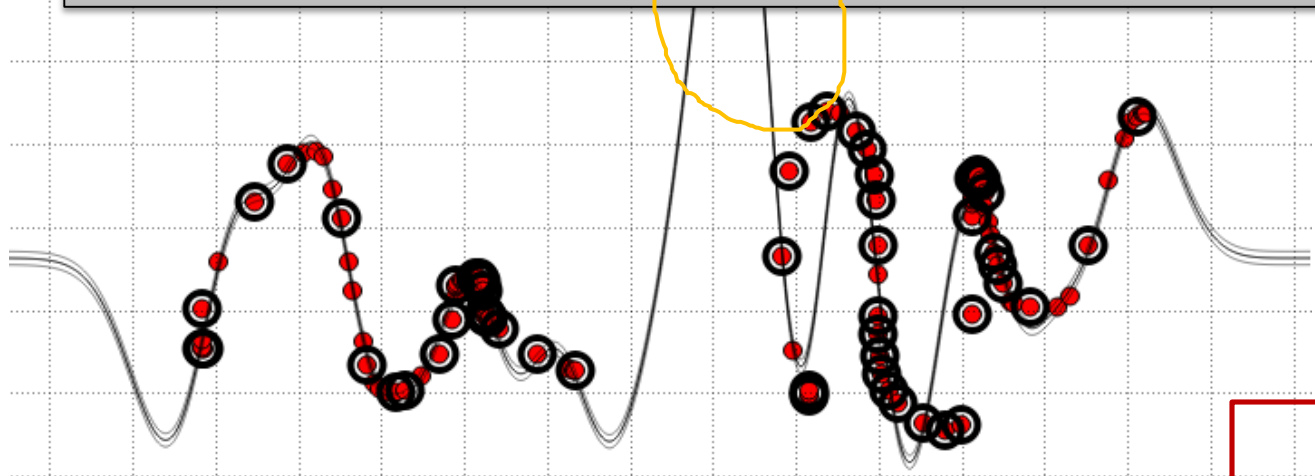
C becomes an upper bound on the absolute value of the  $\alpha$ !



# Fit with a small penalty factor C=0.9



When performing grid search on the open parameters, choose C to fit within the deviation of  $y$  and *not order of magnitude bigger as in the above example!*



C becomes an upper bound on the absolute value of the  $\alpha$ !

$$\alpha_i \in \left[ 0, \frac{C}{M} \right]$$

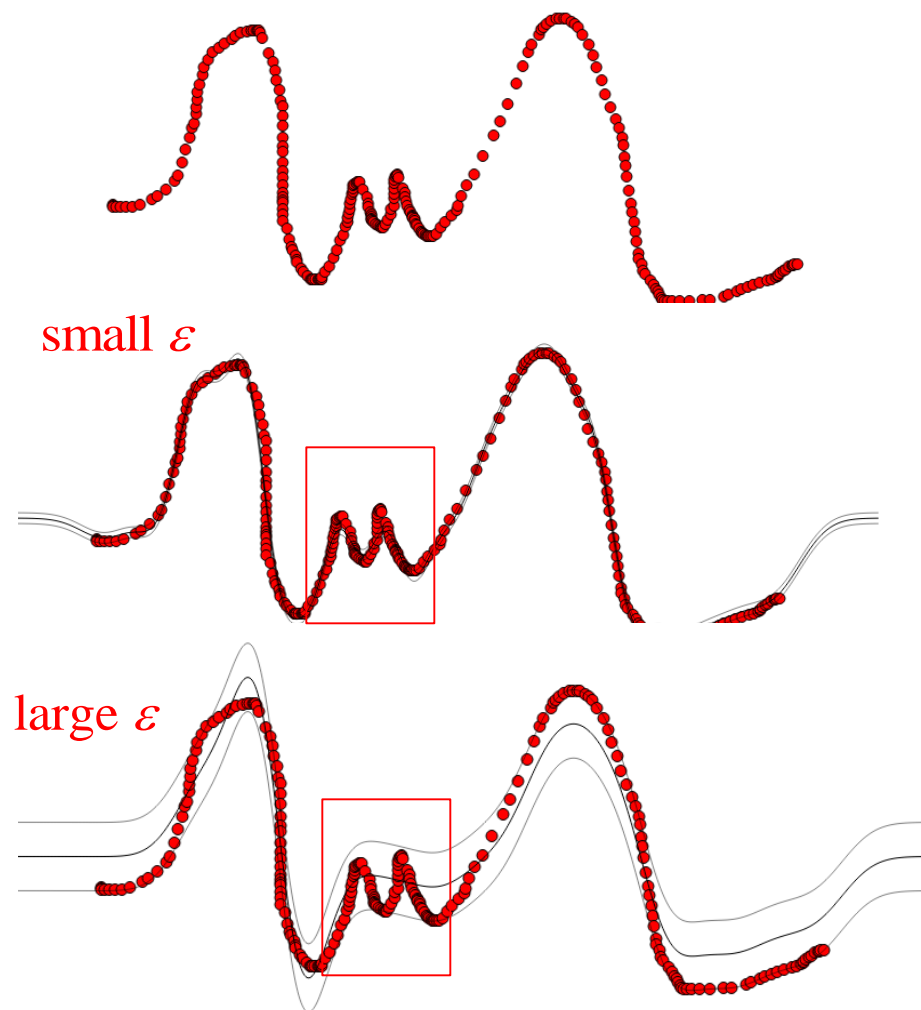
# Fit with a penalty factor C=10



# Recap: SVR Hyperparameters

$\varepsilon > 0$  controls the precision of the fit

$$\begin{aligned} &\text{minimize } \frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{i=1}^M (\xi_i + \xi_i^*) \\ &\text{subject to } \begin{cases} \langle w, \phi(x^i) \rangle + b - y^i \leq \varepsilon + \xi_i \\ y^i - \langle w, \phi(x^i) \rangle - b \leq \varepsilon + \xi_i^* \\ \xi_i \geq 0, \quad \xi_i^* \geq 0 \end{cases} \end{aligned}$$



When performing grid search on the open parameters, choose  $\varepsilon$  to be within the noise observed on  $y$  and the kernel width (for RBF kernels) within the variance on  $x$ .



# Support Vector Regression: $\nu$ -SVR

The standard version of SVR is referred to as  $\varepsilon$ -SVR.

$\nu$ -SVM exploits and rewrites the problem as a convex optimization expression:

$$\min_{w, \xi, \varepsilon} \left( \frac{1}{2} \|w\|^2 + \left[ \underbrace{\nu \varepsilon}_{\text{yellow}} + \frac{1}{M} \sum_{j=1}^M \left( \xi_j + \xi_j^* \right) \right] \right) \quad \text{under constraints}$$

$$\left( w^T \cdot x^j + b \right) - y^j \geq \varepsilon + \xi_j,$$

Penalizes increases of  $\varepsilon$ .

$$y^j - \left( w^T \cdot x^j + b \right) \geq \varepsilon + \xi_j^*,$$

As  $\varepsilon$  decreases, the number of points outside the  $\varepsilon$ -tube increases. The last term penalizes decreases of  $\varepsilon$ .

$$\varepsilon \geq 0, \quad \underline{0 \leq \nu \leq 1}, \quad \xi_j \geq 0, \quad \xi_j^* \geq 0.$$

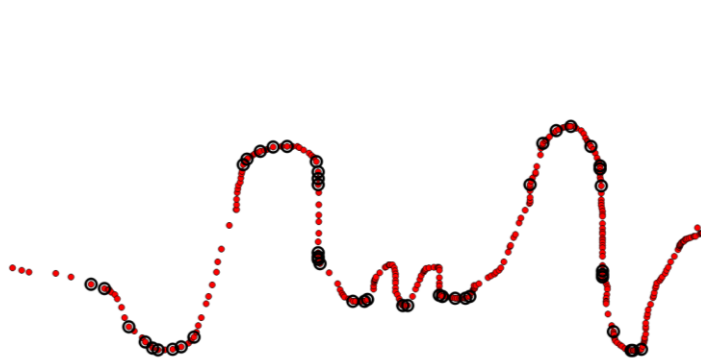
Problem still convex!

$\nu$  is a lower bound on the fraction of support vectors  
(see classification version –  $\nu$ -SVM)



# Support Vector Regression: $\nu$ -SVR

As the number of data grows, so does the number of support vectors.

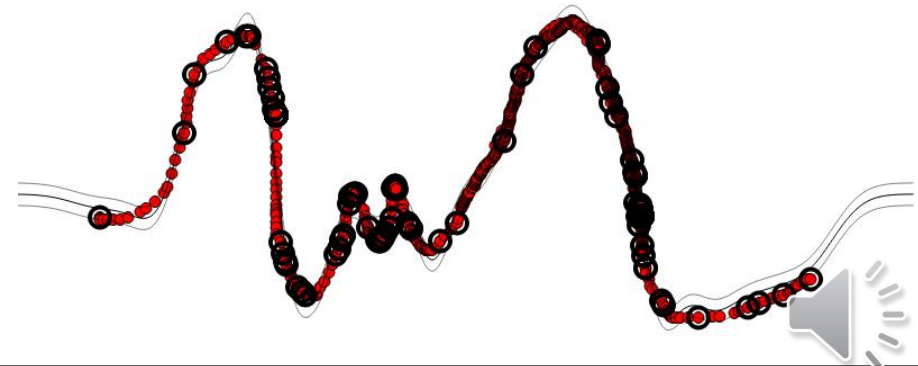


$\nu$ -SVR – 49 SVs

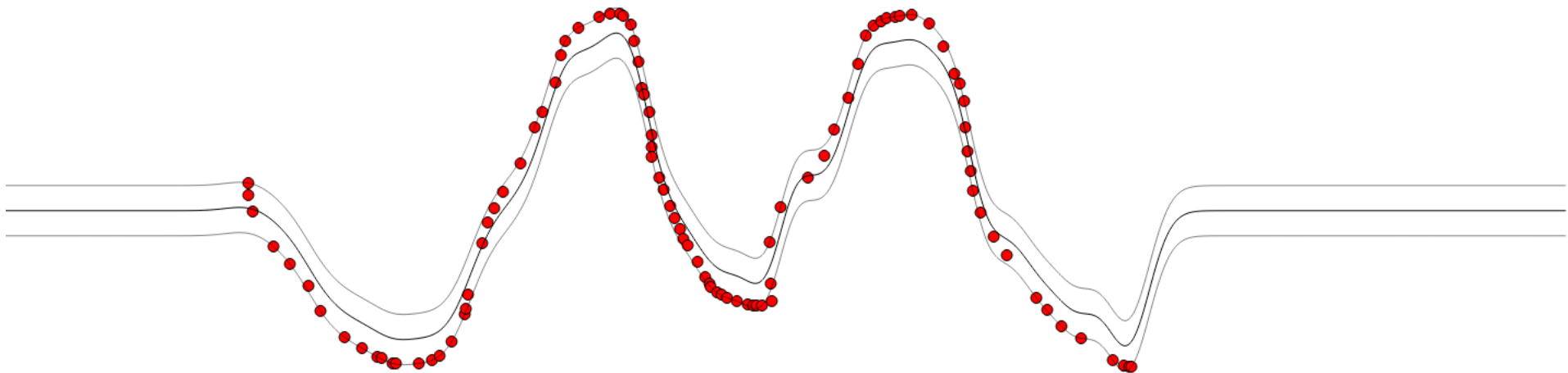
500 datapoints – 79 SVs



1200 datapoints – 117 SVs



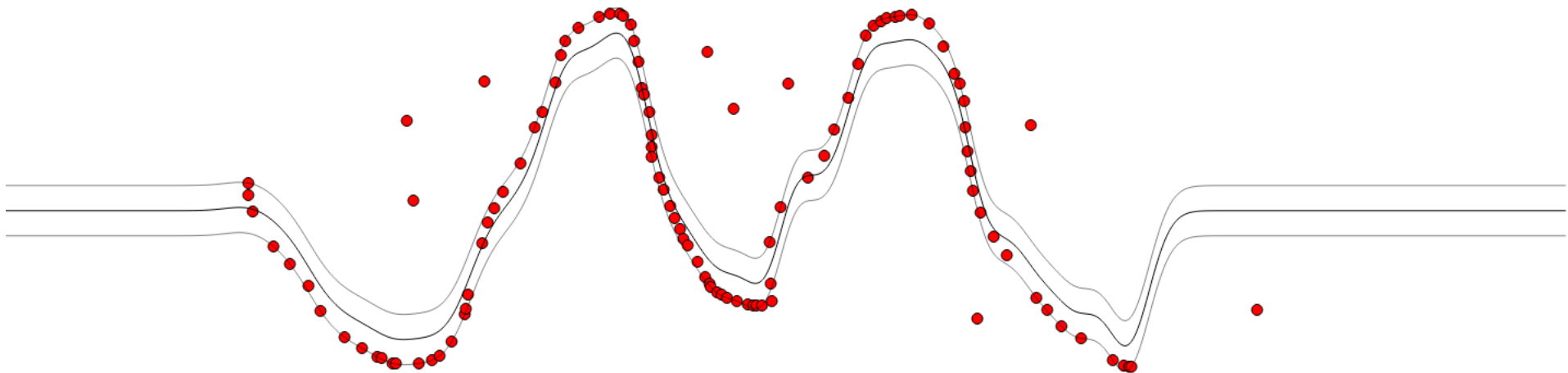
# Support Vector Regression: $\nu$ -SVR



Effect of the automatic adaptation of  $\varepsilon$  using  $\nu$ -SVR



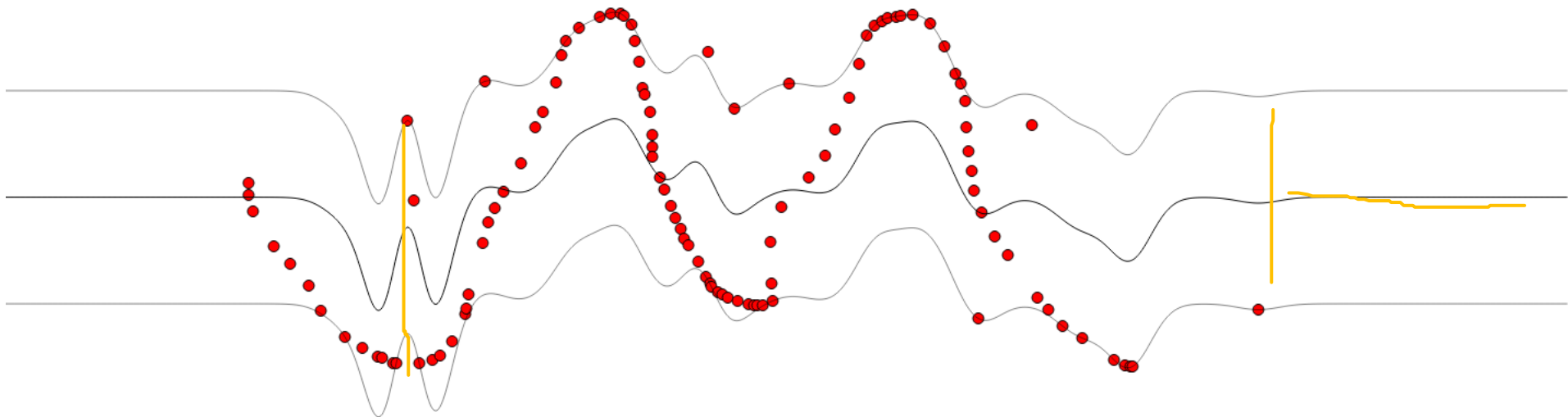
# Support Vector Regression: $\nu$ -SVR



Effect of the automatic adaptation of  $\epsilon$  using  $\nu$ -SVR



# Support Vector Regression: $\nu$ -SVR



Effect of the automatic adaptation of  $\varepsilon$  using  $\nu$ -SVR



# Relevance Vector Regression (RVR)

Relevance Vector Regression (RVR) is the regression version of RVM.

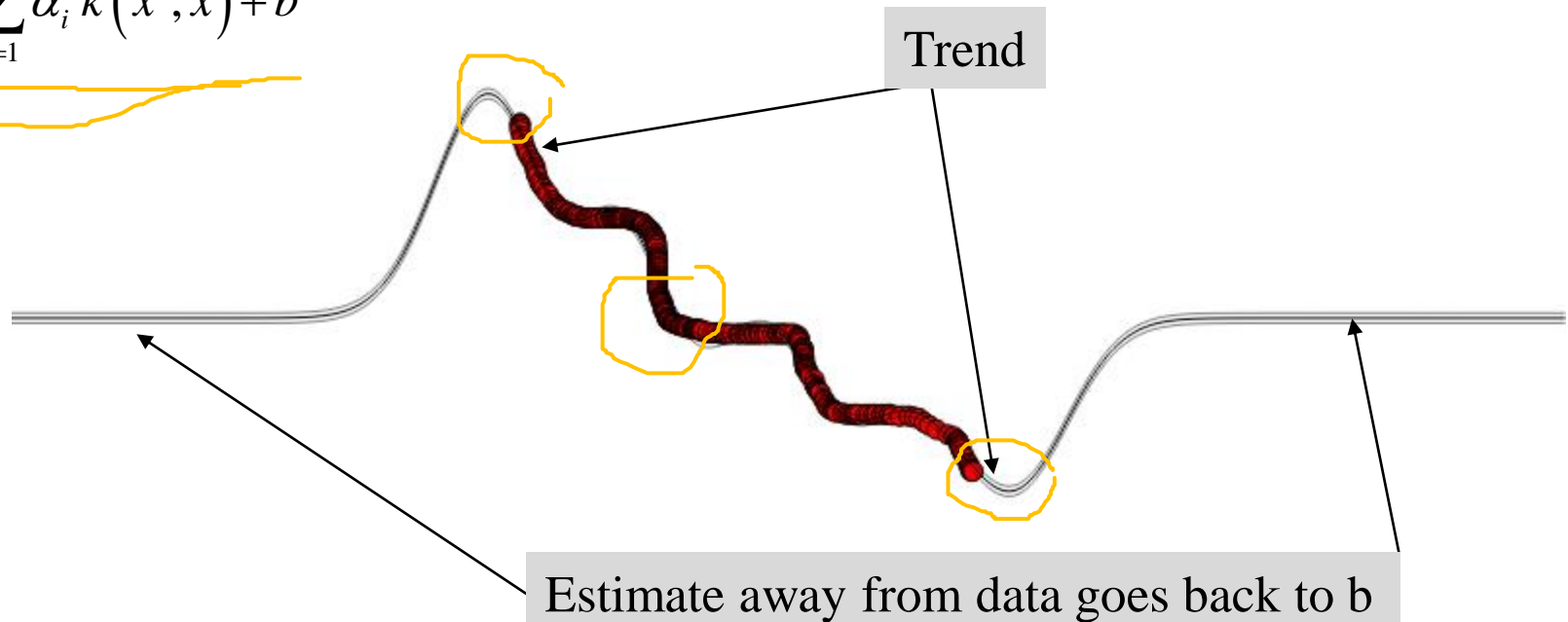
- ❑ It provides a sparser version of SVR.
- ❑ It offers a probabilistic interpretation of the solution.  
(see Tipping 2011, supplementary material to the class)



# $\epsilon$ -SVR - prediction

It is a deterministic model. It does not entail a notion of likelihood.

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



# Relevance Vector Regression (RVR)

Same principle as that described for RVM (see slides on SVM and extensions).  
The derivation of the parameters however differ (see Tipping 2011 for details).

To recall, we start from the solution of SVM.

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

Rewrite the solution of SVM as a linear combination over M basis functions

A sparse solution has a majority of entries with alpha zero.

$$\alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \alpha_M \end{bmatrix} = \begin{bmatrix} \alpha_1 \\ 0 \\ \alpha_3 \\ 0 \\ 0 \\ \cdot \\ \alpha_M \end{bmatrix}$$

In the (binary) classification case,  $y \in [0;1]$ .  
In the regression case,  $y \in \mathbb{R}$ .



# Relevance Vector Regression (RVR)

Rewrite the solution of SVR in a compact form such that the problem is linear in the parameters:

$$y(x) = f(x) = \sum_{i=1}^M \alpha_i \underbrace{k(x, x^i)}_{\psi_i(x)} + \underbrace{\alpha_0 \psi_0(x)}_{b}$$

$$y(x) = \alpha^T \Psi(x), \quad \Psi(x) = \left[ \psi_0(x) \ \psi_1(x) \dots \psi_M(x) \right]^T, \quad \psi_0(x) = 1.$$

Open parameters  $\alpha$ .



# Relevance Vector Regression (RVR)

Take a Bayesian approach and assume that all samples  $y_i$  are i.i.d and that they are measurements of the real value  $\alpha^T \Psi(x^i)$  subjected to white noise  $\varepsilon$ , i.e.:

$$y_i(x) = \alpha^T \Psi(x^i) + \varepsilon, \quad \varepsilon \sim N(0, \sigma_\varepsilon)$$

Likelihood of the model under the dataset  $\{X, \mathbf{y}\}$  composed of

$M$  pairs of data points:  $\{X, \mathbf{y}\} = \{x^i, y^i\}_{i=1}^M$

$$\mathbf{y} = \alpha^T X + N(0, \sigma_\varepsilon)$$

$$\Rightarrow \mathbf{y} \sim p(\mathbf{y} | X, \alpha, \sigma_\varepsilon)$$



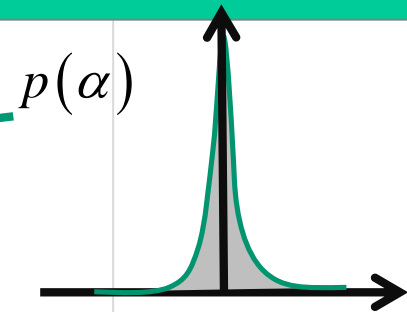
# Relevance Vector Regression (RVR)

Estimates the **posterior** distribution on  $\alpha$ , given the data using Bayes' rule:  $p(\alpha, \sigma_\varepsilon | \mathbf{y}, X)$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$$

$$\Rightarrow p(\alpha, \sigma_\varepsilon | \mathbf{y}, X) = \frac{p(\mathbf{y} | X, \alpha, \sigma_\varepsilon) p(\alpha, \sigma_\varepsilon)}{p(\mathbf{y} | X)}$$

Sparsity is obtained when the distribution is sharply peaked at zero.

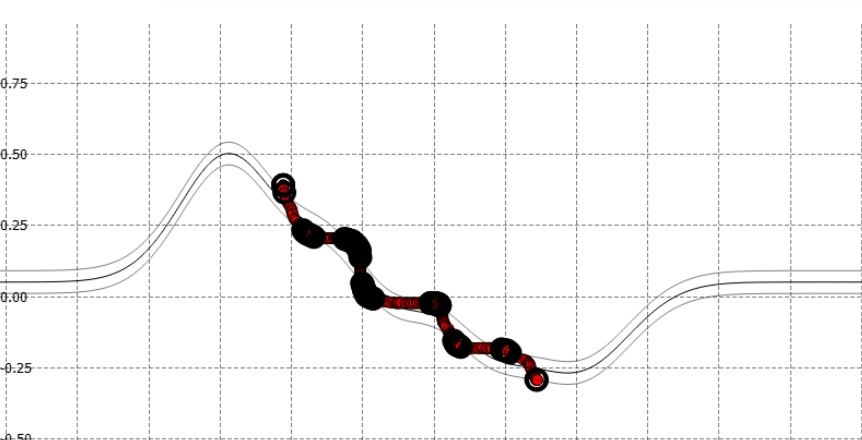


Unlike Ridge Regression Process: No closed-form solution, but iterative procedure.

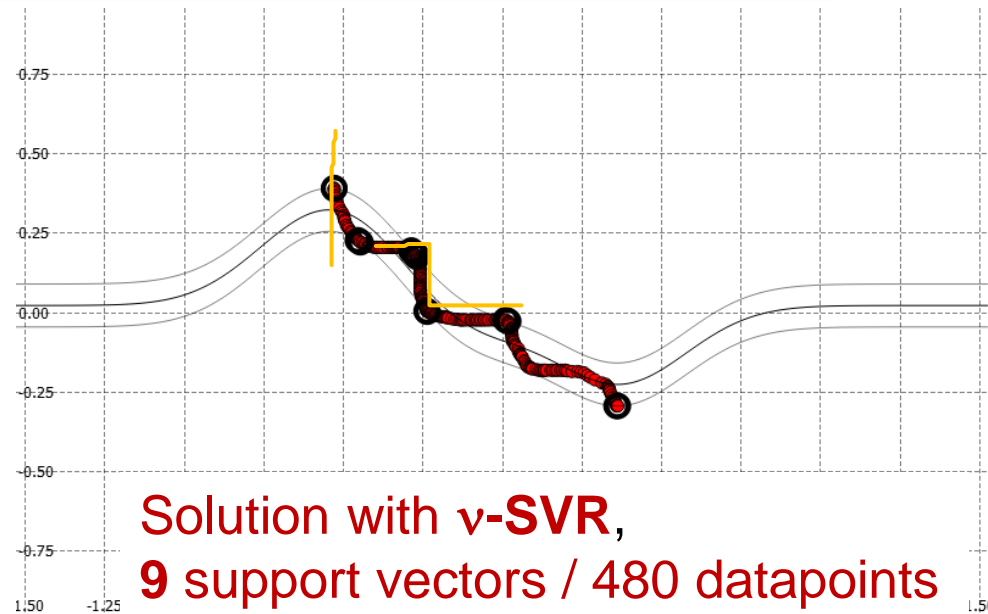
Unlike SVR, not convex optimization any more!



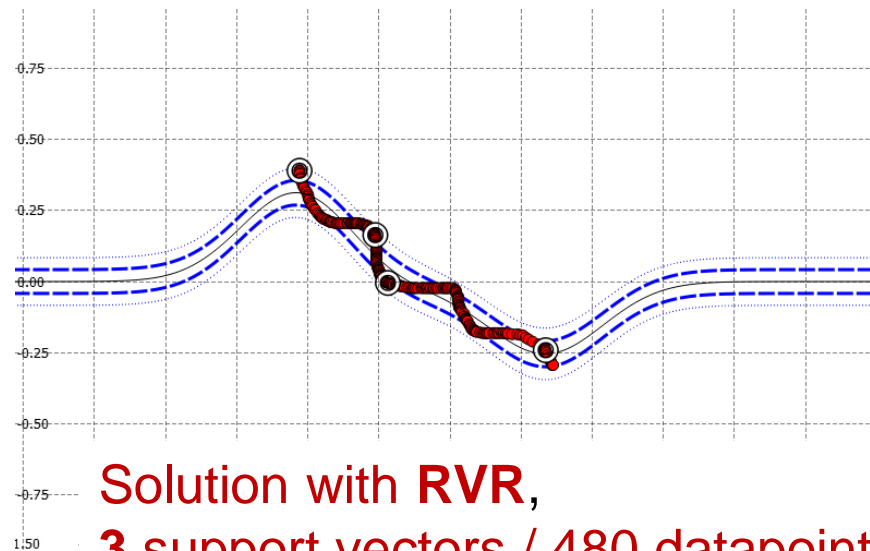
# Comparison $\varepsilon$ -SVR, $\nu$ -SVR, RVR: RBF kernel



Solution with  $\varepsilon$ -SVR,  
87 support vectors / 480 datapoints



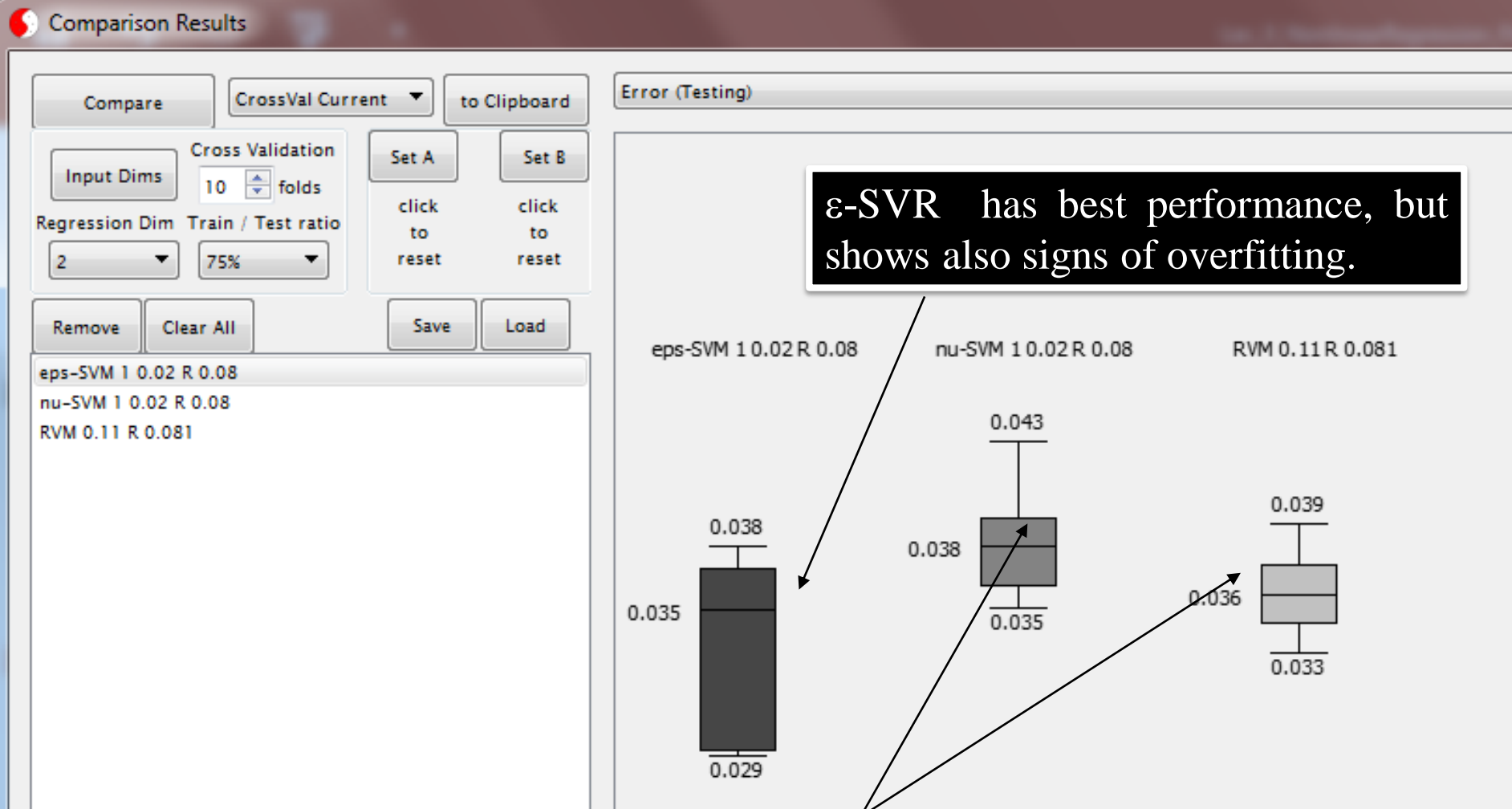
Solution with  $\nu$ -SVR,  
9 support vectors / 480 datapoints



Solution with RVR,  
3 support vectors / 480 datapoints



# Comparison $\epsilon$ -SVR, $\nu$ -SVR, RVR: RBF kernel



RVR has the least number of SVs and better performance than  $\nu$ -SVR.



# Limitations and extensions of $\varepsilon$ -SVR

## Main limitations:

- ❖ Difficult to determine hyperparameters
- ❖ Lack of constraints to limit number of support vectors
- ❖ Lack of probabilistic prediction

## Extensions:

### ☐ $\nu$ -SVR:

- Removes choice of noise model.
- Reduces number of support vectors.

### ❖ RVR:

- Reduces number of support vectors.
- Measure of the goodness of the fit through likelihood.

