

# \*\*\*\*\* Advanced Machine Learning \*\*\*\*\*

## Advanced Regression Techniques

M. Sc. Daniel Wehner

SAP SE / DHBW Mannheim

Summer term 2020

# Agenda for this Unit

<b>Bayesian Regression</b>	<b>4</b>
Introduction	4
 <b>Kernel Ridge Regression</b>	 <b>5</b>
Introduction	5
Woodbury Matrix Identity	6
Example	9
 <b>Gaussian Process Regression</b>	 <b>10</b>
Introduction	10
Learning a Gaussian Process Model	12

Example	15
Learning the Hyper-Parameters	20
<b>Support Vector Regression</b>	<b>21</b>
Introduction	21
Optimization	25
Karush-Kuhn-Tucker Conditions	29
Example	31

# Bayesian Regression

## Introduction

-

# Kernel Ridge Regression

## Introduction

- In ridge regression, the optimal parameters  $\theta$  can be found using the **normal equation**:

$$\theta = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \quad (1)$$

- In the above formula,  $\Phi$  denotes the design matrix (regressor matrix),  $y$  is the label vector and  $\lambda$  is the regularization parameter.
- In order to apply kernels, we have to rephrase this equation in terms of dot products of the input features. Replacing these dot products by kernels avoids operating in feature space.
- This can be achieved by using the **Woodbury matrix identity**.

## Woodbury Matrix Identity

- For the prediction  $y_q$  of a new query data point  $x_q$ , we have to calculate:

$$y_q = \varphi(x_q)^\top \theta \quad (2)$$

Step ①: Insert normal equation  $\Rightarrow$  eq. (1):

$$= \varphi(x_q)^\top (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top y \quad (3)$$

Step ②: Apply Woodbury matrix identity:

$$= \varphi(x_q)^\top \Phi^\top (\Phi \Phi^\top + \lambda I)^{-1} y \quad (4)$$

- The formula given in  $\Rightarrow$  eq. (4) exclusively uses dot products of input features and is therefore susceptible to kernels.

- Replace the dot products by kernel functions:

Rewrite of  $\varphi(\mathbf{x}_q)^\top \Phi^\top$ :

$$\varphi(\mathbf{x}_q)^\top \Phi^\top = \varphi(\mathbf{x}_q)^\top \begin{bmatrix} \varphi(\mathbf{x}^{(1)})^\top \\ \vdots \\ \varphi(\mathbf{x}^{(n)})^\top \end{bmatrix}^\top = \begin{bmatrix} \mathcal{K}(\mathbf{x}_q, \mathbf{x}^{(1)}) \\ \vdots \\ \mathcal{K}(\mathbf{x}_q, \mathbf{x}^{(n)}) \end{bmatrix} = \mathbf{K}_*(\mathbf{x}_q) \quad (5)$$

Rewrite of  $\Phi \Phi^\top$ :

$$\Phi \Phi^\top = \begin{bmatrix} \varphi(\mathbf{x}^{(1)})^\top \\ \vdots \\ \varphi(\mathbf{x}^{(n)})^\top \end{bmatrix} \begin{bmatrix} \varphi(\mathbf{x}^{(1)})^\top \\ \vdots \\ \varphi(\mathbf{x}^{(n)})^\top \end{bmatrix}^\top = \begin{bmatrix} \mathcal{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & \mathcal{K}(\mathbf{x}^{(n)}, \mathbf{x}^{(1)}) \\ \vdots & \ddots & \vdots \\ \mathcal{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(n)}) & \dots & \mathcal{K}(\mathbf{x}^{(n)}, \mathbf{x}^{(n)}) \end{bmatrix} = \mathbf{K} \quad (6)$$

- The kernel matrices  $\mathbf{K}$  and  $\mathbf{K}_*$  must fulfill **Mercer's condition** and therefore have to be **positive-semi definite (psd)**. Famous choices: Polynomial kernel or radial basis function (RBF) kernel.

- The final kernel ridge regression formula is given by:

$$y_q = K_*(x_q)(K + \lambda I)^{-1}y \quad (7)$$

- Like all kernel methods, it is a **non-parametric** approach.



**Kernel methods do not work well for very large data sets ( $> 10,000$  data points), since we have to calculate all pairwise similarities!**



## Example

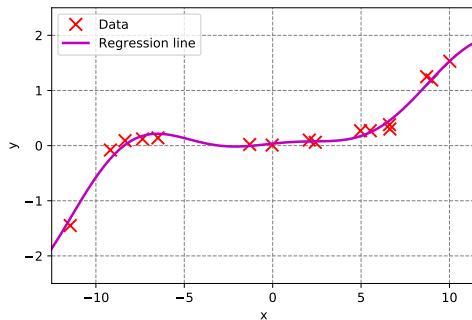


Figure 1:

Result of kernel ridge regression

# Gaussian Process Regression

## Introduction

- Similarly to kernel ridge regression, Gaussian processes do not make any assumptions about the type of regression function (e. g. linear, quadratic, ...)
- It is non-parametric and a form of supervised learning:

$$h(\mathbf{x}) = \mathcal{GP}(m(\mathbf{x}), \mathcal{K}(\mathbf{x}, \mathbf{x}')) \quad (8)$$

- In  $\Rightarrow$  eq. (8),  $m(\mathbf{x})$  denotes the mean function, whereas  $\mathcal{K}(\mathbf{x}, \mathbf{x}')$  denotes the kernel function, which – in the context of Gaussian processes – is referred to as the covariance function.
- Definition of a Gaussian process:  
*Formally, a Gaussian process is a collection of random variables, any finite number of which has a **joint Gaussian distribution**.*

- Instead of modeling a distribution over parameters (cf. Bayesian regression), we model a **distribution over possible regression functions**.
- Thus, Gaussian processes extend multivariate Gaussian distributions to **infinite dimensions**.
  - E. g. a function  $f : \mathbb{R} \mapsto \mathbb{R}$  can be thought of as a sample from some infinite Gaussian distribution.
  - Pick the function which maximizes the posterior distribution over functions.
- The mean of the prior  $m(\mathbf{x})$  distribution is usually set to 0 everywhere.
- In practice, the squared exponential function ( $\hat{=}$  RBF-kernel) is frequently used:

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \cdot \exp \left\{ \frac{-\|\mathbf{x} - \mathbf{x}'\|^2}{2 \cdot l^2} \right\} \quad (9)$$

- Hyper-Parameters:
  - $\sigma_f^2$  denotes the maximum allowable covariance. It should be high for functions covering a broad range of the  $y$ -axis. If  $\mathbf{x} \approx \mathbf{x}'$ ,  $\mathcal{K}(\mathbf{x}, \mathbf{x}')$  approaches this maximum.
  - $l$  (landmark) controls how much the data points influence each other.

## Learning a Gaussian Process Model

- We are given a training data set  $\mathcal{D}$  comprising  $n$  observations:

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$$

- Also, we have a query data point  $\mathbf{x}_q$ , for which  $y_q$  has to be predicted.
- To do so, we compute the covariance between all example pairs.
- This results in three matrices  $\mathbf{K}$  (matrix),  $\mathbf{K}_*$  (vector) and  $\mathbf{K}_{**}$  (scalar).

The matrices have the following form:

$$K = \begin{bmatrix} \mathcal{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \mathcal{K}(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & \dots & \mathcal{K}(\mathbf{x}^{(n)}, \mathbf{x}^{(1)}) \\ \mathcal{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & \mathcal{K}(\mathbf{x}^{(2)}, \mathbf{x}^{(2)}) & \dots & \mathcal{K}(\mathbf{x}^{(n)}, \mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(n)}) & \mathcal{K}(\mathbf{x}^{(2)}, \mathbf{x}^{(n)}) & \dots & \mathcal{K}(\mathbf{x}^{(n)}, \mathbf{x}^{(n)}) \end{bmatrix} \quad (10)$$

$$K_* = [\mathcal{K}(\mathbf{x}_q, \mathbf{x}^{(1)}) \quad \mathcal{K}(\mathbf{x}_q, \mathbf{x}^{(2)}) \quad \dots \quad \mathcal{K}(\mathbf{x}_q, \mathbf{x}^{(n)})]^\top \quad (11)$$

$$K_{**} = \mathcal{K}(\mathbf{x}_q, \mathbf{x}_q) \quad (12)$$



**$K$  is a matrix (contains the similarities of training data pairs),  $K_*$  is a vector (contains similarities of the query data point with the training data), while  $K_{**}$  is actually a scalar (comparison of data point  $\mathbf{x}_q$  to itself)!**

- Since we assume that the data can be modeled as a sample from a multivariate Gaussian distribution, we can model the Gaussian process prior as follows:

$$\begin{bmatrix} \mathbf{y} \\ y_q \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_*^\top \\ \mathbf{K}_* & K_{**} \end{bmatrix} \right) \quad (13)$$

- What we actually want is the **posterior distribution**  $p(y_q | \mathbf{y})$ : 'Given the data, what is  $y_q$ ?'
  - For Gaussian distributions, the posterior distribution can be computed analytically:

$$y_q | \mathbf{y} \sim \mathcal{N} \left( \underbrace{\mathbf{K}_* \mathbf{K}^{-1} \mathbf{y}}_{\text{Matrix of regr. coeff.}}, \underbrace{K_{**} - \mathbf{K}_* \mathbf{K}^{-1} \mathbf{K}_*^\top}_{\text{Schur complement}} \right) \quad (14)$$

- The mean of the posterior distribution is given by the **matrix of regression coefficients**, its variance can be computed using the **Schur complement**.
- We can compute confidence intervals (e. g. 90 % | 95 % | 99 %):

$$(1.65 \mid 1.96 \mid 2.58) \cdot \sqrt{\text{var}(y_q)} \quad (15)$$

## Example

x	y
-1.50	-1.60
-0.25	0.50
0.00	0.80
1.00	-2.00
5.00	0.00
5.50	1.00
10.50	3.00
11.50	3.00

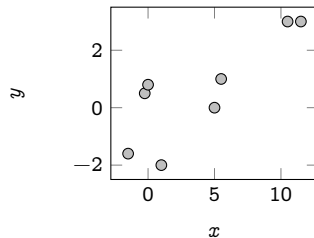
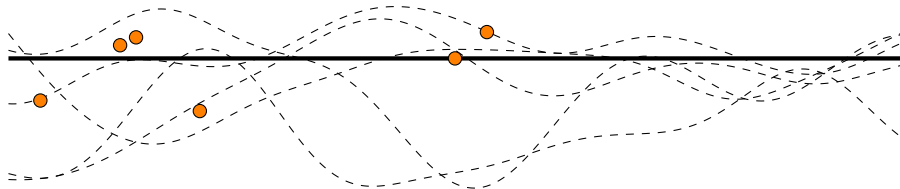


Figure 2:

Example data set for a Gaussian process

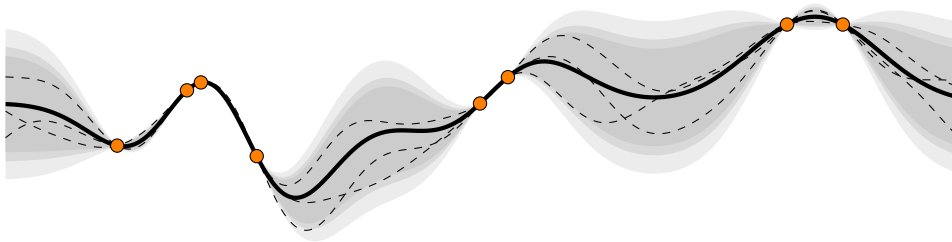
- Suppose  $\sigma_f = 1.27, l = 1.00$ . What is  $y_q$  for  $x_q = 8$ ?
- Let's plot the prior distribution first.

**Prior distribution****Figure 3:**

Prior distribution for the Gaussian process

- Naturally, the prior does not fit the data well (we have not fitted the model yet).
- We have zero mean everywhere.



**Posterior distribution****Figure 4:**

Posterior distribution for the Gaussian process

**Wait a minute: Isn't this model overfitting the training data?**

- The model clearly overfits the data as can be seen from the previous slide (the regression line goes through each training data point perfectly).
- This is because the model assumes the data to be **noise-free**.
- It is possible to add a little bit of noise, in order to deal with this easily ( $\sigma_n$  is the variance of the noise):

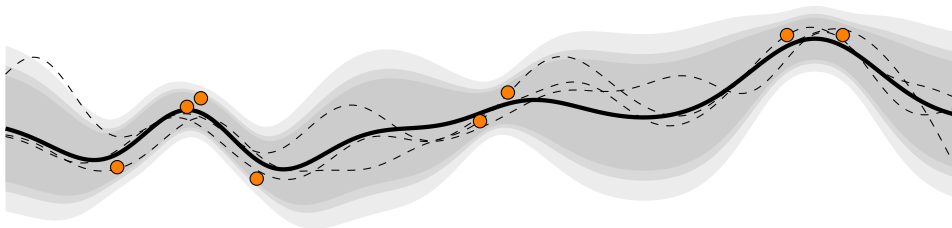
$$K_{\sigma_n} \longleftarrow K + \sigma_n I \quad (16)$$

- The updated formulas look like this:
  - Matrix of regression coefficients (**same result as in kernel ridge regression**):

$$K_* K_{\sigma_n}^{-1} y \quad (17)$$

- Schur complement:

$$K_{**} - K_* K_{\sigma_n}^{-1} K_*^\top \quad (18)$$

**Prior distribution (with noise)****Figure 5:**

Posterior distribution for the Gaussian process with noise

## Learning the Hyper-Parameters

- The results of Gaussian process regression depend heavily on the parameters  $\{\sigma_f, l\}$ , which is why these parameters should be optimized for the task at hand.
- This can be done by maximizing the **marginal likelihood** (e. g. by using gradient ascent).



The exact procedure is very involved and out of scope for this lecture.

# Support Vector Regression

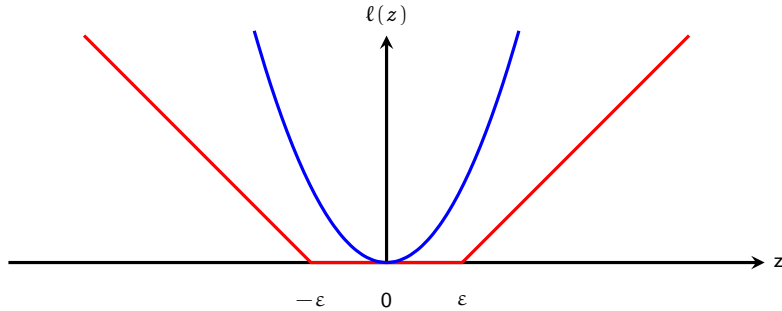
## Introduction

- Support vector machines can be extended to regression problems, while preserving the property of sparseness.
- In ordinary least squares, we minimize a regularized error function given by:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (\hat{h}(\mathbf{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (19)$$

- In the following,  $\boldsymbol{\theta} = \{\mathbf{w}, b\}$  and  $\hat{h}(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\varphi}(\mathbf{x}) + b$ .
- To obtain sparse solutions, the quadratic error is replaced by an  **$\varepsilon$ -insensitive error function**, which gives zero error if the absolute difference between the prediction and the target is less than  $\varepsilon$ :

$$\ell_\varepsilon(\hat{h}(\mathbf{x}) - y) = \begin{cases} 0 & \text{if } |\hat{h}(\mathbf{x}) - y| < \varepsilon \\ |\hat{h}(\mathbf{x}) - y| - \varepsilon & \text{otherwise} \end{cases} \quad (20)$$

**Figure 6:**

An  $\varepsilon$ -insensitive error function (red) compared to the quadratic error function (blue)

- We therefore minimize a regularized error function given by:

$$\mathcal{J}(\boldsymbol{\theta}) = C \sum_{i=1}^n \ell_{\varepsilon}(\hat{h}(\mathbf{x}^{(i)}) - y^{(i)}) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (21)$$

- Analogously to support vector machines for classification,  $C$  denotes the (inverse) regularization parameter.
- Again, we introduce slack variables:
  - We now need two slack variables  $\xi_i \geq 0$  and  $\hat{\xi}_i \geq 0$  for each data point  $\mathbf{x}^{(i)}$ .
  - $\xi_i > 0$  corresponds to a point for which  $y^{(i)} > \hat{h}(\mathbf{x}^{(i)}) + \varepsilon$ .
  - $\hat{\xi}_i \geq 0$  corresponds to a point for which  $y^{(i)} < h(\mathbf{x}^{(i)}) - \varepsilon$ .
- The error function for support vector regression can then be rewritten as:

$$\mathcal{J}(\boldsymbol{\theta}) = C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (22)$$

Illustration of SVM regression, showing the regression curve together with the  $\varepsilon$ -insensitive 'tube'. Also shown are examples of the slack variables  $\xi$  and  $\hat{\xi}$ .

Points above the  $\varepsilon$ -tube have  $\xi > 0$  and  $\hat{\xi} = 0$ , points below the tube have  $\xi = 0$  and  $\hat{\xi} > 0$ . Points inside the tube are characterized by  $\xi = \hat{\xi} = 0$ .

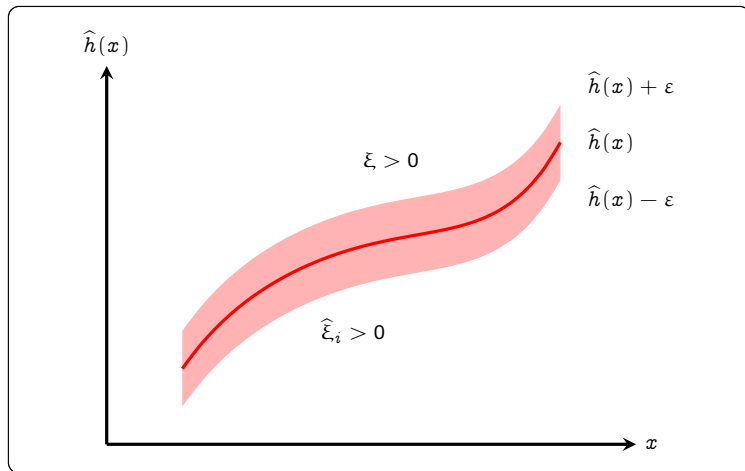
**Figure 7:**

Illustration of support vector regression



## Optimization

- The cost function given by  $\Rightarrow$  eq. (22) must be minimized subject to the constraints:

$$\xi_i \geq 0 \quad (23)$$

$$\hat{\xi}_i \geq 0 \quad (24)$$

$$y^{(i)} \leq h(\mathbf{x}^{(i)}) + \varepsilon + \xi_i \quad (25)$$

$$y^{(i)} \geq h(\mathbf{x}^{(i)}) - \varepsilon - \hat{\xi}_i \quad (26)$$

- This can be achieved by introducing Lagrange multipliers  $\alpha_i \geq 0$ ,  $\hat{\alpha}_i \geq 0$ ,  $\mu_i \geq 0$  and  $\hat{\mu}_i \geq 0$ :

$$\begin{aligned} \mathcal{L} = & C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n (\mu_i \xi_i + \hat{\mu}_i \hat{\xi}_i) \\ & - \sum_{i=1}^n \alpha_i (\varepsilon + \xi_i + \hat{h}(\mathbf{x}^{(i)}) - y^{(i)}) - \sum_{i=1}^n \hat{\alpha}_i (\varepsilon + \hat{\xi}_i - \hat{h}(\mathbf{x}^{(i)}) + y^{(i)}) \end{aligned} \quad (27)$$

**Derivatives of  $\mathcal{L}$** 

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \stackrel{!}{=} 0 \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \varphi(\mathbf{x}^{(i)}) \quad (28)$$

$$\frac{\partial \mathcal{L}}{\partial b} \stackrel{!}{=} 0 \quad \Rightarrow \quad \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) = 0 \quad (29)$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} \stackrel{!}{=} 0 \quad \Rightarrow \quad \alpha_i + \mu_i = C \quad (30)$$

$$\frac{\partial \mathcal{L}}{\partial \hat{\xi}_i} \stackrel{!}{=} 0 \quad \Rightarrow \quad \hat{\alpha}_i + \hat{\mu}_i = C \quad (31)$$



We can use these results to obtain the dual formulation which has to be maximized.

**Dual formulation**

- The dual formulation is given by:

$$\mathcal{L}(\alpha, \hat{\alpha}) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) \mathcal{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) - \varepsilon \sum_{i=1}^n (\alpha_i + \hat{\alpha}_i) + \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) y^{(i)} \quad (32)$$

- The dual is expressed in terms of a kernel function  $\mathcal{K}(\mathbf{x}, \mathbf{x}')$ .
- Maximize the dual function:  $\max_{\alpha, \hat{\alpha}} \mathcal{L}(\alpha, \hat{\alpha})$
- Again, this is a constraint optimization problem which is optimized subject to:

$$0 \leq \alpha_i \leq C \quad (33)$$

$$0 \leq \hat{\alpha}_i \leq C \quad (34)$$

- We again have the **box constraints** which directly follow from the fact that the Lagrange multipliers have to be  $\geq 0$  together with  $\Rightarrow$  eq. (30) and  $\Rightarrow$  eq. (31).

- Substituting  $\Rightarrow$  eq. (28) into  $\hat{h}(\mathbf{x})$ , we see that predictions for new inputs can be made using:

$$\hat{h}(\mathbf{x}) = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \mathcal{K}(\mathbf{x}, \mathbf{x}^{(i)}) + b \quad (35)$$

- The support vectors are those data points for which  $\alpha_i \neq 0$  or  $\hat{\alpha}_i \neq 0$ . Such points either lie on the boundary of the  $\varepsilon$ -tube or outside the tube. All points within the tube have  $\alpha_i = \hat{\alpha}_i = 0$ .
- It is again a **sparse solution**, since we only need the support vectors for the prediction.

## Karush-Kuhn-Tucker Conditions

- The **Karush-Kuhn-Tucker (KKT) conditions** state that at the solution, the product of dual variables and constraints must vanish.
- The KKT conditions for support vector regression are given by:

$$\alpha_i (\varepsilon + \xi_i + \widehat{h}(\mathbf{x}^{(i)}) - y^{(i)}) = 0 \quad (36)$$

$$\widehat{\alpha}_i (\varepsilon + \widehat{\xi}_i - \widehat{h}(\mathbf{x}^{(i)}) + y^{(i)}) = 0 \quad (37)$$

$$\overbrace{(C - \alpha_i)}^{\mu_i} \xi_i = 0 \quad (38)$$

$$\underbrace{(C - \widehat{\alpha}_i)}_{\widehat{\mu}_i} \widehat{\xi}_i = 0 \quad (39)$$



We can derive useful results from the KKT conditions (cf. next slide).

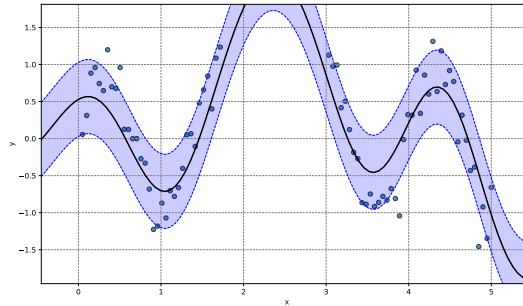
- First of all, we note that  $\alpha_i$  can only be **non-zero**, if  $\varepsilon + \xi_i + \widehat{h}(\mathbf{x}^{(i)}) - y^{(i)} = 0$ . This implies that the data point either lies on the upper boundary of the  $\varepsilon$ -tube ( $\xi_i = 0$ ) or above it ( $\xi_i > 0$ ).
- Analogous:  $\widehat{\alpha}_i$
- The two constraints  $\varepsilon + \xi_i + \widehat{h}(\mathbf{x}^{(i)}) - y^{(i)}$  and  $\varepsilon + \widehat{\xi}_i - \widehat{h}(\mathbf{x}^{(i)}) + y^{(i)}$  are incompatible. This can be seen by adding them together and noting that  $\xi_i, \widehat{\xi}_i$  are non-negative and  $\varepsilon$  is strictly positive. So for every data point  $\mathbf{x}^{(i)}$ , either  $\alpha_i$  or  $\widehat{\alpha}_i$  (or both) must be zero.
- Parameter  $b$  in  $\Rightarrow$  eq. (35) can be found by considering a data point for which  $0 < \alpha_i < C$  ( $\hat{=}$  support vector). From  $\Rightarrow$  eq. (38) it must have  $\xi_i = 0$ . Therefore, according to  $\Rightarrow$  eq. (36) it must satisfy  $\varepsilon + \widehat{h}(\mathbf{x}^{(i)}) - y^{(i)} = 0$ .
- For  $b$  we obtain:

$$b = y^{(i)} - \varepsilon - \mathbf{w}^\top \varphi(\mathbf{x}^{(i)}) \quad (40)$$

$$= y^{(i)} - \varepsilon - \sum_{j=1}^n (\alpha_j - \widehat{\alpha}_j) \mathcal{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \quad (41)$$

- In practice, it is better to consider all support vectors to find  $b$  (average).

## Example



**Figure 8:**

Example of support vector regression using `scikit-learn`