



UNIVERSIDADE
CATÓLICA
PORTUGUESA

BRAGA

Machine Learning

Session 19 - T

Support Vector Machines – Part 3

Ciência de Dados Aplicada

2023/2024

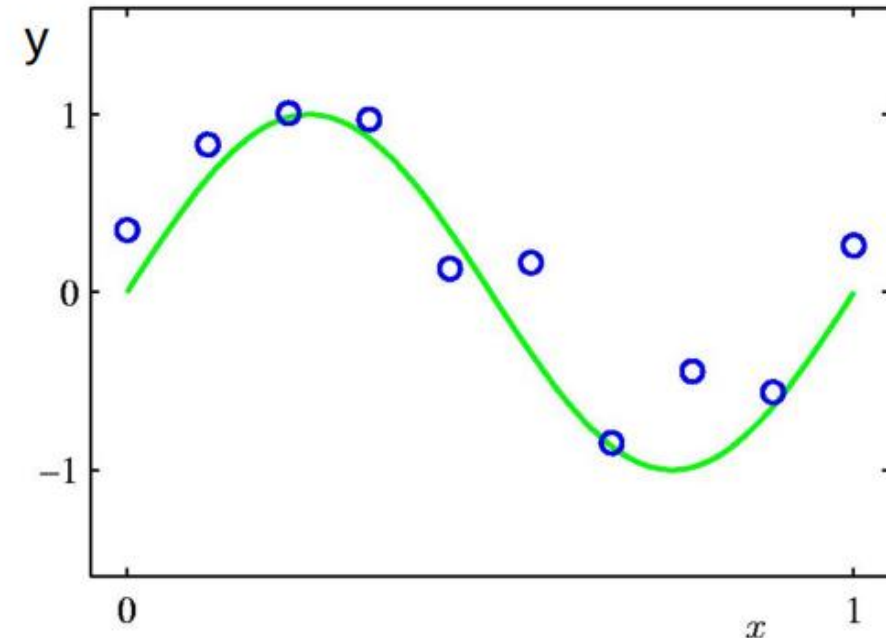
SVMs - Regression

- Suppose we are given a training set of N observations

$$((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)) \text{ with } \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

- The **regression** problem is to estimate $f(x)$ from the data such that

$$y_i = f(\mathbf{x}_i)$$



SVMs - Regression

- As for classification, learning a regressor can be formulated as an optimization problem:

- Minimize
$$\sum_{i=1}^N \underbrace{l(f(\mathbf{x}_i), y_i)}_{\text{loss function}} + \underbrace{\lambda R(f)}_{\text{regularization}}$$

- There is a choice of both **loss function** and **regularization**
 - e.g. squared loss, "Hinge-like" loss
 - ridge, lasso regularization

SVMs - Choice of Regression Function

- Function for regression $y(x, \mathbf{w})$ is a non-linear function of x , but linear in \mathbf{w} :

$$f(\mathbf{x}, \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + w_2\phi_2(\mathbf{x}) + \dots + w_M\phi_M(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x})$$

- For example, for $x \in \mathbb{R}$, polynomial regression with $\phi_j(x) = x^j$:

$$f(x, \mathbf{w}) = w_0 + w_1\phi_1(x) + w_2\phi_2(x) + \dots + w_M\phi_M(x) = \sum_{j=0}^M w_j x^j$$

e.g. for $M = 3$,

$$f(x, \mathbf{w}) = (w_0, w_1, w_2, w_3) \begin{pmatrix} 1 \\ x \\ x^2 \\ x^3 \end{pmatrix} = \mathbf{w}^\top \Phi(x)$$

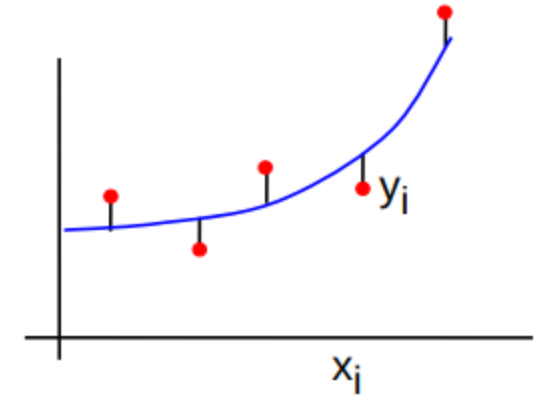
$$\Phi : x \rightarrow \Phi(x) \quad \mathbb{R}^1 \rightarrow \mathbb{R}^4$$

SVMs - Least Squares Ridge Regression

- Cost function – squared loss:

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \underbrace{\{f(x_i, \mathbf{w}) - y_i\}^2}_{\text{loss function}} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{regularization}}$$

target value



- Regression function for x:

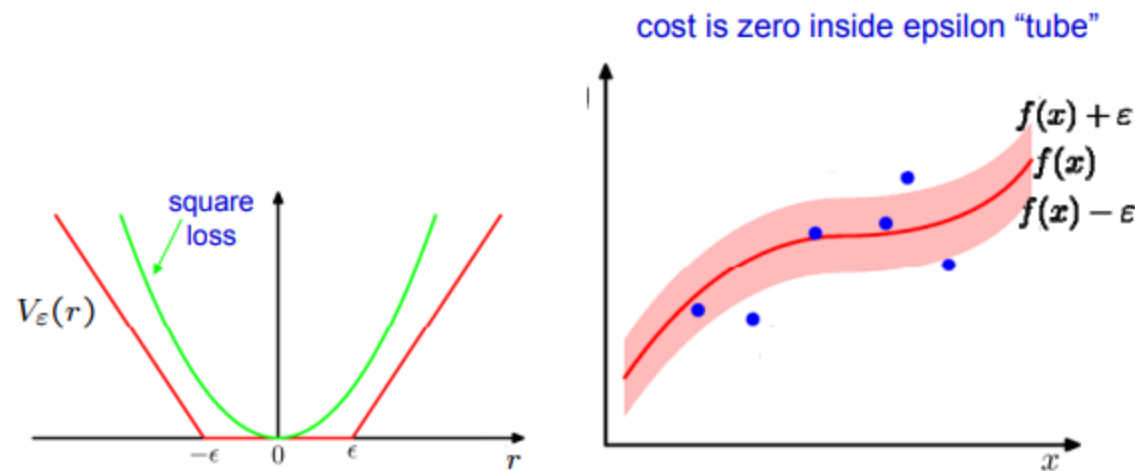
$$f(\mathbf{x}, \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + w_2 \phi_2(\mathbf{x}) + \dots + w_M \phi_M(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x})$$

SVMs - Loss Function for Regression

- To allow for misclassification in SVM regression, we can use the **ϵ -insensitive loss**:

$$J_{\epsilon} = \sum_{i=1}^m J_{\epsilon}(\mathbf{x}_i), \text{ where}$$

$$J_{\epsilon}(\mathbf{x}_i) = \begin{cases} 0 & \text{if } |y_i - (\mathbf{w} \cdot \mathbf{x}_i + w_0)| \leq \epsilon \\ |y_i - (\mathbf{w} \cdot \mathbf{x}_i + w_0)| - \epsilon & \text{otherwise} \end{cases}$$



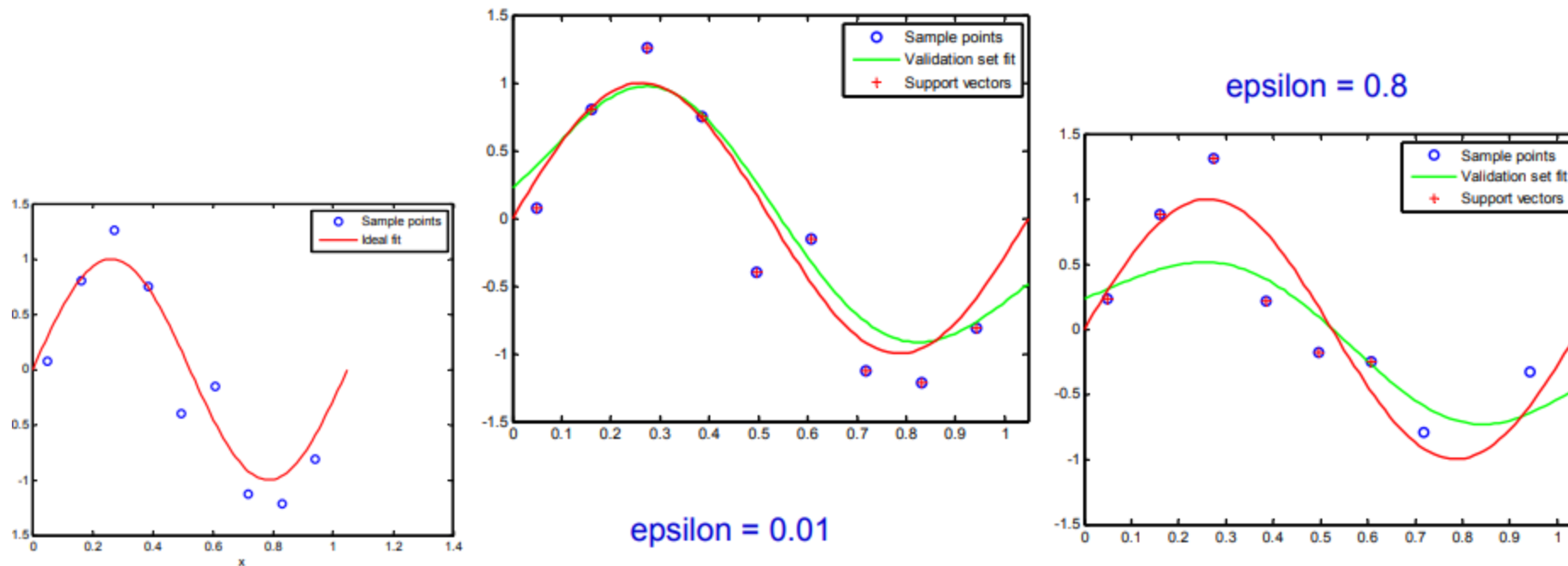
SVMs - The Optimization Problem

$$\begin{array}{ll}\min & \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_i (\xi_i^+ + \xi_i^-) \\ \text{w.r.t.} & \mathbf{w}, w_0, \xi_i^+, \xi_i^- \\ \text{s.t.} & y_i - (\mathbf{w} \cdot \mathbf{x}_i + w_0) \leq \epsilon + \xi_i^+ \\ & y_i - (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq -\epsilon - \xi_i^- \\ & \xi_i^+, \xi_i^- \geq 0\end{array}$$

- As before, Kernels can be used to get non-linear functions

SVMs - Effect of ϵ

- As ϵ increases, the function is allowed to move away from the data points, the number of support vectors decreases and the fit gets worse.



Resources

- <https://www.youtube.com/watch?v=kPw1IGUAoY8>
- https://www.researchgate.net/publication/228537532_Support_Vector_Regression