

Introduction to Machine Learning - CS725  
Instructor: Prof. Ganesh Ramakrishnan  
Lecture 8 - Support Vector Regression and  
Optimization Basics

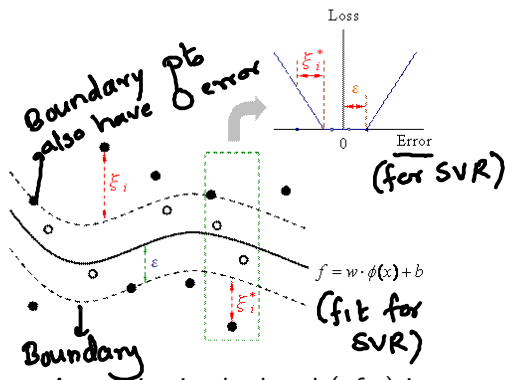
# Building on questions on Least Squares Linear Regression

- 1 Is there a probabilistic interpretation?
  - Gaussian Error, Maximum Likelihood Estimate
- 2 Addressing overfitting
  - Bayesian and Maximum A posteriori Estimates, Regularization, Support Vector Regression
- 3 How to minimize the resultant and more complex error functions?
  - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

# Support Vector Regression

One more formulation before we look at [Tools of Optimization/duality](#)

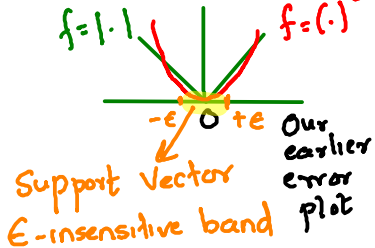
# Support Vector Regression (SVR)



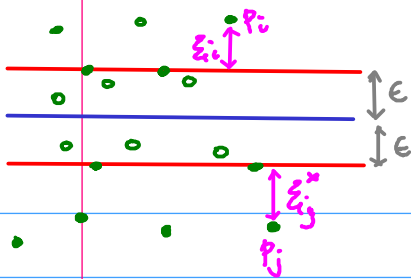
Loss function so far

$$= f(w^T \phi(x_i) - y_i)$$

$$f = 1 \cdot | \cdot |$$
$$f = (\cdot)^2$$



- Any point in the band (of  $\epsilon$ ) is not penalized. Thus the loss function is known as  $\epsilon$ -insensitive loss
- Any point outside the band is penalized, and has slackness  $\xi_i$  or  $\xi_i^*$
- The SVR model curve may not pass through any training point



For all pts with  $\epsilon$  band,  
we expect  $\xi_i = \xi_i^* = 0$   
The regression curve  
may not pass through any  
point at all!

You need both  $\xi_i$   
&  $\xi_i^*$  for each  $p_i$  since  
 $p_i$  could lie above or below

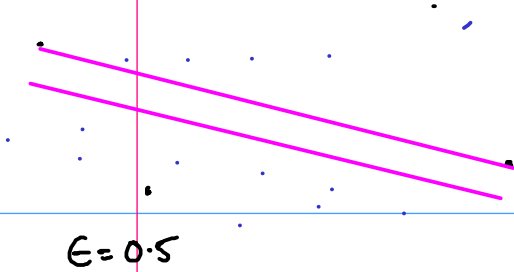
$$p_i \rightarrow \xi_i, \xi_i^*$$

$$p_j \rightarrow \xi_j, \xi_j^*$$

$\xi_i$  = distance of  $p_i$  above the  $\epsilon$  band

$\xi_i^*$  = distance of  $p_i$  below the  $\epsilon$ -band

} A pt can be either  
above OR below  
 $\therefore \xi_i > 0 \Rightarrow \xi_i^* = 0$   
and vice versa



H/w: Think of a setting in which no point will land up in the  $E$ -band

$$f(x) = \omega^T \phi(x) + b$$
$$\xi_i, \xi_i^+, \epsilon, \mathcal{D} = \{(x_i, y_i)\}$$

- The tolerance  $\epsilon$  is fixed
- It is desirable that  $\forall i$ :

Constraints?

Objective desired to be minimized?

## Constraints First!

- The tolerance  $\epsilon$  is fixed
- It is desirable that  $\forall i$ :

$$\textcircled{1} \bullet y_i - w^T \phi(x_i) - b \leq \epsilon + \xi_i$$

$$\bullet b + w^T \phi(x_i) - y_i \leq \epsilon + \xi_i^*$$

$$\textcircled{a} y_i, \xi_i = y_i - w^T \phi(x_i) - b - \epsilon$$

$\updownarrow \epsilon$

$$\textcircled{b} \bullet y_i, y_i - w^T \phi(x_i) - b - \epsilon \leq 0$$

$$\textcircled{a} \bullet y_i, w^T \phi(x_i) + b - y_i - \epsilon \leq 0$$

$$\textcircled{c} y_i, \xi_i^* = w^T \phi(x_i) + b - y_i - \epsilon$$

Even if  $y_i \geq w^T \phi(x_i) + b$   
 [Summarizes cases  $\textcircled{a}$  &  $\textcircled{b}$ ]

Even if  $y_i \leq w^T \phi(x_i) + b$   
 [summarizes cases  $\textcircled{c}$  &  $\textcircled{a}$ ]

Alternatively:

$$\xi_i = \max(y_i - w^T \phi(x_i) - b - \epsilon, 0)$$

$$\xi_i^* = \max(b + w^T \phi(x_i) - y_i - \epsilon, 0)$$



# SVR objective

- 1-norm Error, and  $L_2$  regularized:

$$\frac{1}{2} \|\omega\|_2^2 + C \left( \sum_{i=1}^m \xi_i + \xi_i^* \right)$$

$L_2$  regularizer

Think of  $C$  as  $\frac{1}{2\lambda}$

1-norm Error

$$\text{st } y_i - \omega^T \phi(x_i) - b \leq \epsilon + \xi_i \quad (1)$$

$$\omega^T \phi(x_i) + b - y_i \leq \epsilon + \xi_i^* \quad (2)$$

$$\xi_i, \xi_i^* \geq 0$$

$\sum_{i=1}^m (\xi_i + \xi_i^*)$  = sum of errors over all  $m$  data points

Q1: Why is  $\xi_i, \xi_i^* \geq 0$  required?

Q2: How does this formulation ensure that one of  $\xi_i$  &  $\xi_i^*$  will be 0?

Ans to Q2: Proof by contradiction.

↳ Suppose  $\xi_i > 0$  &  $\xi_i^* > 0$

$$\Leftrightarrow y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i \quad (1)$$

$$w^\top \phi(x_i) + b - y_i \leq \epsilon + \xi_i^* \quad (2)$$

Claim:  $\xi_i^* = 0$  would also satisfy (2) [complete proof]

$\Rightarrow$  By setting  $\xi_i^* = 0$ , I can lower my error in

objective while still satisfying (1) and (2)

$\Rightarrow$  Contradicts assumption that  $\xi_i > 0$   $\xi_i^* > 0$  was an "optimal" solution!

- 1-norm Error, and  $L_2$  regularized:

- $\min_{w, b, \xi_i, \xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i + \xi_i^*)$   
s.t.  $\forall i,$   
 $y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i,$   
 $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*,$   
 $\xi_i, \xi_i^* \geq 0$

- 2-norm Error, and  $L_2$  regularized:

$$+ C \sum_i (\xi_i^2 + \xi_i^{*2})$$

$\xi_i, \xi_i^* \geq 0 \Leftarrow$  Not required! (H/W)

- 1-norm Error, and  $L_2$  regularized:

- $\min_{w,b,\xi_i,\xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i + \xi_i^*)$   
s.t.  $\forall i,$   
 $y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i,$   
 $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*,$   
 $\xi_i, \xi_i^* \geq 0$

- 2-norm Error, and  $L_2$  regularized:

- $\min_{w,b,\xi_i,\xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2})$   
s.t.  $\forall i,$   
 $y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i,$   
 $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*$
- Here, the constraints  $\xi_i, \xi_i^* \geq 0$  are not necessary

# Need for Optimization so far

- Unconstrained (**Penalized**) Optimization:

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|_2^2 + \Omega(\mathbf{w})$$

- **Constrained Optimization 1:**

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|_2^2$$

such that  $\Omega(\mathbf{w}) \leq \theta$

- **Constrained Optimization 2** ( $t = 1$  or  $2$ ):  $\rightarrow$  SVR

$$\arg \min_{w, b, \xi_i, \xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^t + \xi_i^{*t})$$

s.t.  $\forall i, y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i$ ;  $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*$

- **Equivalence:**  $\lambda$  (**Penalized**)  $\equiv$   $\theta$  (**Constrained**)

- **Duality:** Dual of Support Vector Regression: **Kernels, non-linear**

**Non-parametric**

# Solving Unconstrained Minimization Problem

- Intuitively: Minimize by setting derivative (gradient) to 0 and hoping to find **closed form** solution.
- When is such a solution a global minimum?
- For most optimization problems, finding closed form solutions is difficult. Even for linear regression (for which closed form solution exists), are there alternative methods?
  - Eg: Consider,  $\mathbf{y} = \phi\mathbf{w}$ , where  $\phi$  is a matrix with full column rank, the least squares solution,  $\mathbf{w}^* = (\phi^T\phi)^{-1}\phi^T\mathbf{y}$ . Now, imagine that  $\phi$  is a very large matrix. with say, 100,000 columns and 1,000,000 rows. Computation of closed form solution might be challenging.
- How about iterative methods?