

Introduction to Machine Learning - CS725

Instructor: Prof. Ganesh Ramakrishnan

Lecture 07 - 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 for Gaussian and Laplacian (and Beta) priors, L_0 , L_1 and L_2 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

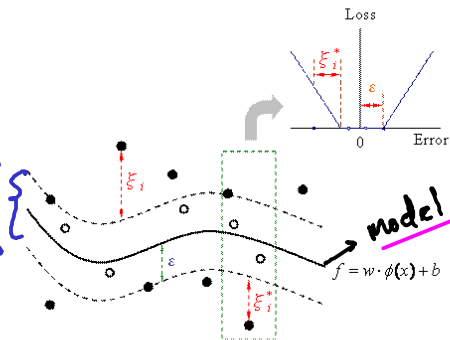
→ SVR motivates Duality (kernelized representation)
→ Equivalences → Penalized regression & constrained regression

Support Vector Regression

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

Support Vector Regression (SVR)

Region of ϵ -error



So far: Generalize (avoid overfitting) through $\Omega(w)$ (L_1, L_2, L_0)

Now: Avoid overfitting (generalize) also through Loss component

$$\min_w \underline{\text{Loss}(w, \mathcal{D}) + \Omega(w)}$$

- Any point in the band (of ϵ) is not penalized. Thus the loss function is known as ϵ -insensitive loss

- Any point outside the band is penalized, and has slackness ξ_i or ξ_i^*

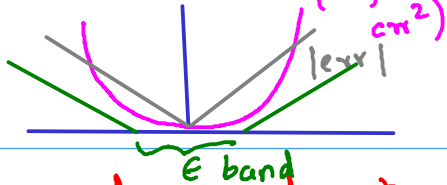
- The SVR model curve may not pass through any training point

↳ Acknowledging that measurement errors be ignored

In words: Give me a regression curve (in x)

$$f(x) = w^T \phi(x) + b$$

such that:

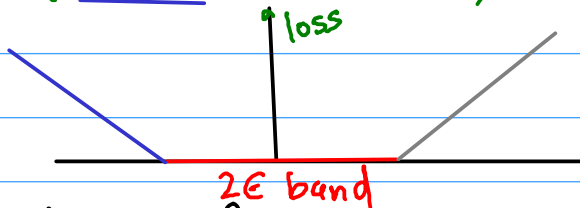


- ① Training pts in an ϵ -band around curve do not contribute to any loss (New part)
- ② Loss/Error contribution from pts outside the ϵ -band is minimized (as $\epsilon \rightarrow 0$ problem tends to earlier problem)
- ③ $\Omega(w)$ is penalty (continued from before)

We will formulate & solve SUR problem, which, as $\epsilon \rightarrow 0$ also helps solve problem of $\min |err|$ discussed in TUTORIAL 1

Form of SVR ϵ errors:

$$\text{loss} = \max(\text{err} - \epsilon, \underline{0}, \text{err} - \epsilon)$$



Constrained form for above

$$\text{loss}(x_i) \geq 0$$

$$\geq -(f(x_i) - y_i) - \epsilon$$

$$\geq (f(x_i) - y_i) - \epsilon$$

} & { min loss(x_i)

$y_i > f(x_i)$ outside ϵ

Intent

$$\xi_i = \max(y_i - f(x_i) - \epsilon, 0) = \max(\text{err}_i - \epsilon, 0)$$

$f(x_i) > y_i$ outside ϵ

$$\xi_i^+ = \max(f(x_i) - y_i - \epsilon, 0) = \max(-\text{err}_i - \epsilon, 0)$$

- The tolerance ϵ is fixed

- It is desirable that $\forall i$:

Constraints

$$\& \xi_i \xi_i^+ = 0$$

pt i can
be either above
OR below ϵ band!

Execution

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

$$0 \leq \xi_i$$

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

$$0 \leq \xi_i^+$$

and

$$\min \sum_i \xi_i + \sum_i \xi_i^+$$

ϵ is a hyperparameter in as much as λ was!

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

- $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i$
- $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$

$$4 \min_{\xi_i, \xi_i^*} \sum_i (\xi_i + \xi_i^*)$$

$(\xi_i \geq 0, \xi_i^* = 0)$
 $(\xi_i = \xi_i^* = 0)$
 $(\xi_i^* = \xi_i = 0)$
 $(\xi_i^* \geq 0, \xi_i \leq 0)$

① $y_i, \xi_i = y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon$

② $y_i, y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon \leq 0$

③ $y_i, \mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i - \epsilon \leq 0$

④ $y_i, \xi_i^* = \mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i - \epsilon$

Prove: That $\xi_i \xi_i^* \leq 0$

Geometric
view of
intent }

Claim: If $\xi_i > 0$ & $\xi_i^* > 0$.. prove contradiction result

$$\left. \begin{array}{l} \textcircled{1} \quad y_i - \overset{k}{\omega^\top \phi(x_i)} - b - \epsilon \leq \xi_i \\ \textcircled{2} \quad \underset{-k}{\omega^\top \phi(x_i) + b - y_i} - \epsilon \leq \xi_i^* \\ 0 \leq \xi_i \\ 0 \leq \xi_i^* \end{array} \right\}$$

$$\min_i (\xi_i + \xi_i^*)$$

$$\rightarrow k - \epsilon = \xi_i \Rightarrow k > \epsilon$$

Proof :- $\xi_i > 0 \Rightarrow \textcircled{1}$ is an equality $\Rightarrow k - \epsilon > 0 \Rightarrow -k - \epsilon < 0$

$\Rightarrow \xi_i^* = 0$ yields lower value of objective while also satisfying $\textcircled{2}$

SVR objective

- 1-norm Error, and L_2 regularized:

$$\min_i C(\underbrace{\sum_i (\xi_i + \xi_i^+)}_{\text{Absolute value of errors outside } \epsilon \text{ band}}) + \underbrace{\|W\|_2^2}_{\text{L}_2 \text{ regularizer}}$$

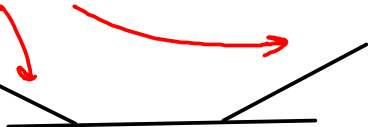
} C is tradeoff between loss & regularizers

$$\text{s.t. } \forall i: y_i - \underbrace{w^T \phi(x_i)}_{-b} - \epsilon \leq \xi_i$$

$$w^T \phi(x_i) + b - y_i$$

$$-\epsilon \leq \xi_i^+$$

$$\xi_i, \xi_i^+ \geq 0$$



SVR objective

- 1-norm Error, and L_2 regularized:

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

- 2-norm Error, and L_2 regularized:

(minimizing)

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2})$$

$$\text{s.t. } y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i$$
$$b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$$

quadratic loss: error^2

Since $\xi_i^2 \rightarrow \infty$ as $\xi_i \rightarrow -\infty$, $\xi_i \geq 0$ is natural. Similarly $\xi_i^* \geq 0$

SVR objective

- 1-norm Error, and L_2 regularized:

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

- 2-norm Error, and L_2 regularized:

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

If ① is satisfied for $\xi_i < 0$
then ① is also satisfied for
 $\xi_i = 0$ while further reducing
the objective ($\sum (\xi_i^2 + \xi_i^{*2})$)

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:**

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

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

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

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

$$\text{s.t. } \forall i, y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i; b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$$

- **Equivalence:** λ (**Penalized**) $\equiv \theta$ (**Constrained**)
- **Duality:** Dual of Support Vector Regression

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 Φ is a matrix with full column rank, the least squares solution, $\mathbf{w}^* = (\Phi^T\Phi)^{-1}\Phi^T\mathbf{y}$. Now, imagine that Φ 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?

- 1-norm Error, and L_2 regularized:

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

- 2-norm Error, and L_2 regularized:

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

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

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

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

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

$$\text{s.t. } \forall i, y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i; b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$$

- **Equivalence:** λ (**Penalized**) $\equiv \theta$ (**Constrained**)
- **Duality:** Dual of Support Vector Regression