# Lecture 09-b: Support Vector Regression in some details

Instructor: Prof. Ganesh Ramakrishnan

Recall motivations for regularizer: $\frac{1}{2}\|w\|_2^2$

① Slides 17 to 19 of https://www.cse.iitb.ac.in/~cs725/notes/lecture-slides/lecture-06-unannotated.pdf

Analytical motivation

Empirical motivation

$\phi w = y$ for $\phi = \begin{bmatrix} 1 & x_i^1 & x_i^2 \cdots x_i^{p-1} \\ 1 & & \\ \vdots & x_m^1 & x_m^2 \cdots x_m^{p-1} \end{bmatrix}$

- will have 1 or more solutions when $p \geq m$ and no solutions (unless rank of $\phi < m$) otherwise. Thus, with $p \geq m$, you have more chances of fitting regression curve on all data points than when $p < m$

- will therefore tend to "overfit" training data at the risk of making errors on test data as $p$ increases

- $p > p' \equiv$ obtaining $w^{p'}$ by setting higher coefficient indices of $w^p$ to zero

$$w^{p'} = w^p(1, 2 \cdots p')$$

- $\therefore$ Determining "right" $p'$ is like finding a "sparse" solution for larger value of $p$

(see pages 3 to 5 of

https://www.cse.iitb.ac.in/~cs725/notes/lecture-slides/lecture-07-annotated.pdf)

Observe how test error increases for $t \geq 7$ whereas train error keeps decreasing for a while.

This indicates that some regularization (to penalize non-zero $w$'s) could have helped avoid overfitting
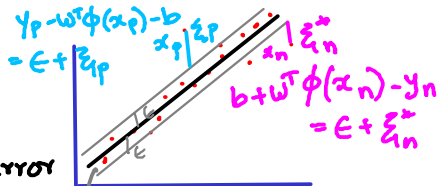
Again see some of the arguments below that are common to the empirical & analytic reasoning

Also see: Slide 12 of https://www.cse.iitb.ac.in/~cs725/notes/lecture-slides/lecture-09-unannotated.pdf
for how the ISTA algo for LASSO attains this sparsity (for $L_1$ regularization)
by thresholding $\hat{w}_i^{(k+1)}$ and only letting "promising" values of $\hat{w}_i^{(k+1)}$ to be
non-zero and setting the rest to zero

# Support Vector (SVR) Regression

$y_p - w^\top \phi(x_p) - b = \epsilon + \xi_p$

$x_q \mid \xi_p^*$

$x_n \mid \xi_n^*$

$b + w^\top \phi(x_n) - y_n = \epsilon + \xi_n^*$

Note: We want

$$\left| \left( y_i - (w^\top \phi(x_i) + b) \right) \right| \leq \epsilon + \xi \text{ error}$$

→ Zero error in this $\epsilon$-band

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

Claim: Exactly one of $\xi_i$ & $\xi_i^* > 0$

Claim: At optimal soln, one of the two will be an equality for all points that lie outside the $\epsilon$-band

- Let's consider the lagrange multipliers $\alpha_i$, $\alpha_i^*$, $\mu_i$ and $\mu_i^*$ corresponding to the above-mentioned constraints respectively.

Claim: For points within $\epsilon$-band, $\xi_i = \xi_i^* = 0$

$$\min_{\omega, b, \xi_i, \xi_i^*} \frac{1}{2}\|\omega\|^2 + C \sum_i^2 (\xi_i + \xi_i^*) \rightarrow ①$$

s.t

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

$$\omega^T \phi(x_i) + b - y_i \leq \epsilon + \xi_i^* \rightarrow ③$$

} for each $i$

$$\xi_i^*, \xi_i \geq 0 \xrightarrow{\hspace{4cm}} ④$$

<u>Claims 1 & 2</u>: Suppose $\xi_i > 0$ then as per ②

$$\underline{y_i - \omega^T \phi(x_i) - b} \leq \epsilon + \xi_i \quad \text{(Multiplying both sides by -1)}$$

$$-y_i + \omega^T \phi(x_i) + b \geq -\epsilon - \xi_i \rightarrow ②'$$

<u>Claim</u>: By setting $\underline{\xi_i^* = 0}$, ③ will be satisfied & objective reduced in contrast with $\xi_i^* > 0$

↳ claim: $y_i - \omega^T \phi(x_i) - b > \epsilon$ (since if LHS $\leq \epsilon$

↳ $-y_i + \omega^T \phi(x_i) + b < -\epsilon < 0$ then $\xi_i = 0$ should

↳ $-y_i + \omega^T \phi(x_i) + b < 0 + \xi_i^*$ have been the soln)

<u>see next slide ↑</u>

Similarly if $\exists \xi_i$ for which $y_i - \omega^T \phi(x_i) - b = \epsilon + \xi_i$ that $\xi_i$ will be chosen

Claim: If $\hat{\xi}_i$ s.t $y_i - \omega^T \phi(x_i) - b < \epsilon + \hat{\xi}_i$

is optimal, then I claim that $\xi_i < \hat{\xi}_i$

$(\xi_i = \epsilon - y_i + \omega^T \phi(x_i) + b)$ with $y_i - \omega^T \phi(x_i) - b = \epsilon + \xi_i$

will give objective $= \frac{1}{2} \|\omega\|_2^2 + C \xi_i + o/rs$

that is less than $\frac{1}{2} \|\omega\|_2^2 + C \hat{\xi}_i + o/rs$

Lagrange fn

$$L(w, b, \xi_i, \xi_i^*, \alpha_i, \alpha_i^*, \mu_i, \mu_i^*) = \frac{1}{2}\|w\|^2 + C\sum(\xi_i + \xi_i^*)$$
$$+ \sum_i \alpha_i(y_i - w^T\phi(x_i) - b - \epsilon - \xi_i) + \sum_i \alpha_i^*(w^T\phi(x_i) + b - y_i - \epsilon - \xi_i^*)$$
$$- \sum_i \mu_i \xi_i - \sum_i \mu_i^* \xi_i^*$$

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

$$s.t \quad y_i - w^T\phi(x_i) - b \le \epsilon + \xi_i \rightarrow \alpha_i \qquad (2)$$

$$w^T\phi(x_i) + b - y_i \le \epsilon + \xi_i^* \rightarrow \alpha_i^* \qquad (3)$$

$$\left. \begin{array}{l} \xi_i \ge 0 \longrightarrow \mu_i \\ \xi_i^* \ge 0 \longrightarrow \mu_i^* \end{array} \right\} (4)$$

$$\nabla_w L(w, b, \xi_i, \xi_i^*, \alpha_i, \alpha_i^*, \mu_i, \mu_i^*) = 0 \Rightarrow w + \sum_i -\alpha_i \phi(x_i) + \alpha_i^* \phi(x_i)$$
$$= 0$$

$$\underset{re}{=} w = \sum_i (\alpha_i^* - \alpha_i)\phi(x_i)$$

$$\nabla_b L(w, b, \xi_i, \xi_i^*, \alpha_i, \alpha_i^*, \mu_i, \mu_i^*) = 0 \Rightarrow \sum_i(\alpha_i - \alpha_i^*) = 0$$

$$\nabla_{\xi_i} L = 0 \Rightarrow C - \alpha_i - \mu_i = 0 \quad \text{similarly: } C - \alpha_i^* - \mu_i^* = 0$$

# KKT conditions

- Differentiating the Lagrangian w.r.t. $w$,
  $w - \alpha_i \phi(x_i) + \alpha_i^* \phi(x_i) = 0$
  i.e. $w = \sum_{i=1}^{n} (\alpha_i - \alpha_i^*) \phi(x_i)$
- Differentiating the Lagrangian w.r.t. $\xi_i$,
  $C - \alpha_i - \mu_i = 0$
  i.e. $\alpha_i + \mu_i = C$
- Differentiating the Lagrangian w.r.t $\xi_i^*$,
  $\alpha_i^* + \mu_i^* = C$
- Differentiating the Lagrangian w.r.t $b$,
  $\sum_i (\alpha_i^* - \alpha_i) = 0$
- Complimentary slackness:
  $\alpha_i(y_i - w^\top \phi(x_i) - b - \epsilon - \xi_i) = 0$
  $\mu_i \xi_i = 0$
  $\alpha_i^*(b + w^\top \phi(x_i) - y_i - \epsilon - \xi_i^*) = 0$
  $\mu_i^* \xi_i^* = 0$

*(handwritten annotations on right side)*

(1) If $\xi_i > 0$ then
$\mu_i = 0$ & $\alpha_i = C$

If $\xi_i^* > 0$ then
$\mu_i^* = 0$ & $\alpha_i^* = C$

(2) $\alpha_i \in (0, C) \Rightarrow$
$\mu_i \in (0, C) \Rightarrow$
$\xi_i = 0$ &
$y_i - w^\top \phi(x_i) - b - \epsilon - \xi_i = 0$
$\Rightarrow y_i - w^\top \phi(x_i) - b = \epsilon$

# Conclusions from the KKT conditions:

$$\alpha_i \in (0, C) \Rightarrow ?$$

$$\alpha_i^* \in (0, C) \Rightarrow ?$$