

Lecture 09: Lasso and Support Vector Regression

Instructor: Prof. Ganesh Ramakrishnan

.

Recap: Duality and KKT conditions

For the previously mentioned formulation of the problem, KKT conditions for all differentiable functions (i.e. f, g_i, h_j) with $\hat{\mathbf{w}}$ primal optimal and $(\hat{\lambda}, \hat{\mu})$ dual optimal point are:

- $\nabla f(\hat{\mathbf{w}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{w}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{w}}) = 0$
- $g_i(\hat{\mathbf{w}}) \leq 0; 1 \leq i \leq m$
- $\hat{\lambda}_i \geq 0; 1 \leq i \leq m$
- $\hat{\lambda}_i g_i(\hat{\mathbf{w}}) = 0; 1 \leq i \leq m$
- $h_j(\hat{\mathbf{w}}) = 0; 1 \leq j \leq p$

Equivalence of the two formulations of regularized least square

Formulation 1:

$$\mathbf{w}^*(\lambda) \xleftarrow{\text{soln}} \min_{\mathbf{w}} f_{\lambda}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \phi \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2 \quad \left. \vphantom{\min_{\mathbf{w}}} \right\} \text{Penalty Version}$$

Formulation 2:

$$\mathbf{w}^*(\eta) \xleftarrow{\text{soln}} \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \phi \mathbf{w}\|^2 \text{ s.t. } \|\mathbf{w}\|^2 - \eta \leq 0 \quad \left. \vphantom{\min_{\mathbf{w}}} \right\} \text{Constrained Version}$$

Handwritten notes: "is soln for some λ " (blue arrow from $\mathbf{w}^*(\eta)$ to $\mathbf{w}^*(\lambda)$), "is soln for some η " (green arrow from $\mathbf{w}^*(\lambda)$ to $\mathbf{w}^*(\eta)$)

The Lagrangian for Formulation 2 is:

$$L(\mathbf{w}, \lambda) = \frac{1}{2} \|\mathbf{y} - \phi \mathbf{w}\|^2 + \alpha (\|\mathbf{w}\|^2 - \eta)$$

Necessary conditions for optimality for Formulation 1 are:

$$\nabla_{\mathbf{w}} f_{\lambda}(\mathbf{w}^*(\lambda)) = 0$$

where $\mathbf{w}^*(\lambda)$ is the optimal solution for a given λ .

Handwritten note: "NO constraint Normal gradient test" (green arrow pointing to the optimality condition)

Equivalence of the two formulations of regularized least square

(constrained)

- For Formulation 2, the KKT conditions imply that we have:
 $\nabla_{\mathbf{w}} f_{\alpha}(\mathbf{w}^*) = \mathbf{0}$ and $\alpha^*(\|\mathbf{w}^*\|^2 - \eta) = \mathbf{0}$, $\alpha^* \geq 0$. (α is Lagrange variable)
- If formulation 1 is solved for a given λ and its solution is $\mathbf{w}^*(\lambda)$ then
 - ▶ by setting $\eta = \|\mathbf{w}^*(\lambda)\|_2^2$, you get that $\alpha^* = \lambda$ and $\mathbf{w}^* = \mathbf{w}^*(\lambda)$ satisfy the KKT conditions for formulation 2, showing that both formulations have the same solution.
 - ▶ if you solved formulation 2 and set $\lambda = \alpha^*$, you attain the same solution as attained by solving Problem 1.



Choice of regularizer and regularization parameter

penalized version

- How do we decide which value of λ to choose for the regularizer? How about choosing the regularization parameter λ through cross-validation?
- Recall the polynomial curve fitting problem we considered earlier. As we increased the degree of the polynomial how will the training error vary? What about the test error? And what is the effect of varying λ on train and test errors? [Extreme case: $\lambda=0 \Rightarrow$ No constraint or penalty \Rightarrow least train error]
- How about a different regularizer?
 - ▶ **Lasso:** When the L_1 norm is used (instead of L_2 as in ridge regression). Replacing $\|w\|_2^2$
- How about a different error function?
 - ▶ **Support Vector Regression.** Replacing $\|\phi w - y\|_2^2$

Cross Validation for choosing λ

$\lambda \in \{\lambda_1, \dots, \lambda_k\}$
Undirected or exhaustive search

Train

Find RMS & choose λ_i with least RMS

Valid-ation

Unknown test data

70%

30%

Alternative:
Binary/directed search:

$$\lambda \in [\lambda_l \quad \lambda_m \quad \lambda_h]$$

$$\text{if } \text{RMS}(\lambda_h) < \text{RMS}(\lambda_m)$$

then now try

$$\lambda \in \left[\lambda_m \quad \frac{\lambda_m + \lambda_h}{2} \quad \lambda_h \right]$$

Lasso: Continuing from Quiz 1, Problem 3



$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \|\phi \mathbf{w} - \mathbf{y}\|^2 \text{ s.t. } \|\mathbf{w}\|_1 \leq \eta, \quad (1)$$

where

$$\|\mathbf{w}\|_1 = \left(\sum_{i=1}^n |w_i| \right) \quad (2)$$

- Since $\|\mathbf{w}\|_1$ is not differentiable, one can express (2) as a set of constraints

$$\sum_{i=1}^n \xi_i \leq \eta, \quad \underline{w_i \leq \xi_i}, \quad \underline{-w_i \leq \xi_i} \quad \left. \begin{array}{l} \text{Supposed that} \\ |w_i| = \xi_i \end{array} \right\}$$

- The resulting problem is a linearly constrained Quadratic optimization problem (LCQP): *(standard optimization solvers exist for LCQP)* $\rightarrow \xi_i: \text{CVX}$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \|\phi \mathbf{w} - \mathbf{y}\|^2 \text{ s.t. } \sum_{i=1}^n \xi_i \leq \eta, \quad \mathbf{w}_i \leq \xi_i, \quad -\mathbf{w}_i \leq \xi_i \quad (3)$$

Common objective fn: $\|\phi w - y\|_2^2$

Constraint 1: $\{w \mid \underbrace{\sum |w_i|}_{C_1} \leq \eta, w \in \mathbb{R}^m\}$

Constraint 2: $\{w \mid \underbrace{w_i \leq \xi_i, -w_i \leq \xi_i, \sum \xi_i}_{C_2} \leq \eta, w \in \mathbb{R}^m\}$

Is $C_1 = C_2$ [i.e. $C_1 \subseteq C_2$ & $C_2 \subseteq C_1$]: Ans Yes!

Lasso: Continued

- KKT conditions:

$\theta_i \rightarrow$ Lagrange mult for $w_i \leq \xi_i$
 $\lambda_i \rightarrow$ Lagrange mult for $-w_i \leq -\xi_i$
 $\beta \rightarrow$ Lagrange mult for $\sum_{i=1}^n \xi_i \leq \eta$

$\nabla L = 0$

$$\leftarrow 2(\phi^T \phi) \mathbf{w} - 2\phi^T \mathbf{y} + \sum_{i=1}^n (\theta_i - \lambda_i) = 0$$

Complementary slackness

$$\leftarrow \beta \left(\sum_{i=1}^n \xi_i - \eta \right) = 0$$

Also: $\theta_i \geq 0$
 $\lambda_i \geq 0$ & $\beta \geq 0$

$$\forall i, \theta_i (w_i - \xi_i) = 0 \text{ and } \lambda_i (-w_i - \xi_i) = 0$$

- Like Ridge Regression, an equivalent Lasso formulation can be shown to be:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \|\phi \mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1 \quad (4)$$

constrained

penalized

- The justification for the equivalence between (2) and (4) as well as the solution to (4) requires *subgradient*.

Subgradients: Generalization of gradient

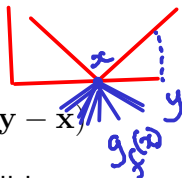


- An equivalent condition for convexity of $f(x)$:

$$\forall \mathbf{x}, \mathbf{y} \in \text{dmn}(f), \mathbf{f}(\mathbf{y}) \geq \mathbf{f}(\mathbf{x}) + \nabla^T \mathbf{f}(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

- $\mathbf{g}_f(\mathbf{x})$ is a *subgradient* for a function f at \mathbf{x} if

$$\forall \mathbf{y} \in \text{dmn}(f), \mathbf{f}(\mathbf{y}) \geq \mathbf{f}(\mathbf{x}) + \mathbf{g}_f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$$



- Any convex (even non-differentiable) function will have a subgradient at any point in the domain!
- If a convex function f is differentiable at \mathbf{x} then $\nabla f(\mathbf{x}) = \mathbf{g}_f(\mathbf{x})$
- \mathbf{x} is a point of minimum of (convex) f if and only if $\mathbf{0}$ is a subgradient of f at \mathbf{x} $\exists \mathbf{g}_f(\mathbf{x}) = \mathbf{0}$

Subgradients and Lasso

- Claim (out of syllabus): If $\mathbf{w}^*(\eta)$ is solution to (2) and $\mathbf{w}^*(\lambda)$ is solution to (4) then
 - ▶ Solution to (2) with $\eta = \|\mathbf{w}^*(\lambda)\|_1$ is also $\mathbf{w}^*(\lambda)$ and
 - ▶ Solution to (4) with λ as solution to $\phi^T(\phi\mathbf{w} - \mathbf{y}) = \lambda \mathbf{g}_x$ is also $\mathbf{w}^*(\eta)$
- The unconstrained form for Lasso in (4) has no closed form solution : $\min_{\mathbf{w}} \|\phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$
- But it can be solved using a generalization of gradient descent called proximal subgradient descent¹

Presented in Elements of Stats ML
(called ISM..)

¹<https://www.cse.iitb.ac.in/~cs709/notes/enotes/lecture27b.pdf>

Proximal Subgradient Descent for Lasso²

(ISTA = Iterative Soft Thresholding in Tibshirani's book/paper)

- Let $\varepsilon(w) = \|\phi w - y\|_2^2$

- Proximal Subgradient Descent Algorithm:**

Initialization: Find starting point $w^{(0)}$

- Let $\hat{w}_u^{(k+1)}$ be a next gradient descent iterate for $\varepsilon(w^{(k)})$

- Compute $w_p^{(k+1)} = \underset{w}{\operatorname{argmin}} \|\underbrace{w - \hat{w}_u^{(k+1)}}_{\text{Proximity}}\|_2^2 + \lambda t \underbrace{\|w\|_1}_{\text{penalty}}$ by

setting subgradient of this objective to 0. This results in:

- If $\hat{w}_i^{(k+1)} > \lambda t$, then $w_i^{(k+1)} = -\lambda t + \hat{w}_i^{(k+1)}$
- If $\hat{w}_i^{(k+1)} < -\lambda t$, then $w_i^{(k+1)} = \lambda t + \hat{w}_i^{(k+1)}$
- 0 otherwise.

Set $k = k + 1$, **until** stopping criterion is satisfied (such as no significant changes in w^k w.r.t $w^{(k-1)}$)

First 2 steps like grad descent

$\hat{w}_u^{(k+1)}$ is unpenalized grad descent update

Solving for $w_p^{(k+1)}$ EXACTLY by setting $g_f = 0$

$-\lambda t$ λt : Feature selection in this soft thresholding

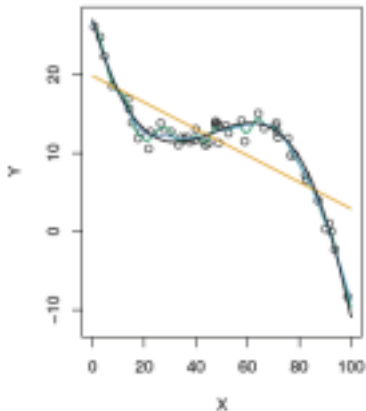
²<https://www.cse.iitb.ac.in/~cs709/notes/enotes/lecture27b.pdf>

Choice of regularizer and regularization parameter

- How do we decide which value of λ to choose for the regularizer? How about choosing the regularization parameter λ through cross-validation?
- Recall the polynomial curve fitting problem we considered earlier. As we increased the degree of the polynomial how will the training error vary? What about the test error? And what is the effect of varying λ on train and test errors?
- How about a different regularizer?
 - ▶ **Lasso:** When the L_1 norm is used (instead of L_2 as in ridge regression).
- How about a different error function?
 - ▶ **Support Vector Regression.**

Support Vector Regression

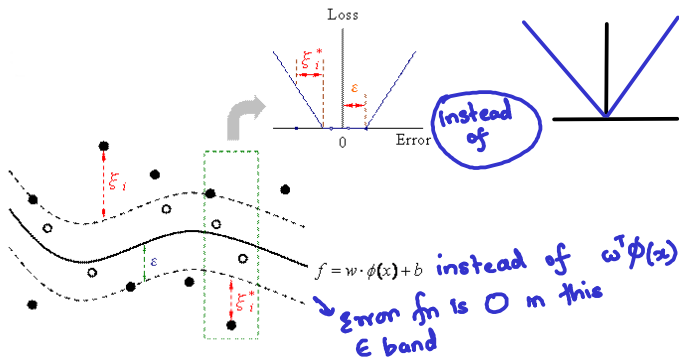
Polynomial regression: $\phi = [x, x^2, x^2, x^2, x, x_2, \dots]$



- Consider a degree 3 polynomial regression model as shown in the figure
- Each bend in the curve
 - corresponds to increase in $\|w\|$
- Eigen values of $(\phi^T \phi + \lambda I)$ are indicative of curvature. Increasing λ reduces the curvature

came from $\lambda \|w\|^2$

Can I get similar through benefit loss fn?



- 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

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

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

ξ_i & ξ_i^* = slack variables

→ deviation when $(y_i - w^T \phi(x_i) - b) > \epsilon$
 } we wanted $|y_i - w^T \phi(x_i) - b| \leq \epsilon$
 ↓ deviation when $(y_i - w^T \phi(x_i) - b) < -\epsilon$

SVR objective

- 1-norm regularized:

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

- 2-norm regularized:

$$\begin{aligned} & \text{min}_{w,b,\xi_i,\xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2}) \\ & \text{s.t. } \forall i, \end{aligned}$$

$$\begin{aligned} & y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i, \\ & b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^* \end{aligned}$$

- ▶ Here, the constraints $\xi_i, \xi_i^* \geq 0$ are not necessary

} Sum of slackness
Variable squared

Try deriving the KKT conditions for the two norm regularized Support Vector Regression problem on slide 18