

Introduction to Machine Learning - CS725

Instructor: Prof. Ganesh Ramakrishnan

Lecture 11 - KKT Conditions, Support Vector Regression and its Dual

KKT conditions for the Constrained (Convex) Problem

- Here, $\mathbf{w} \in \mathbb{R}^n$ and the domain is the intersection of all functions. Lagrangian is:

$$L(\mathbf{w}, \lambda, \mu) = f(\mathbf{w}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{w}) + \sum_{j=1}^p \mu_j h_j(\mathbf{w})$$

- KKT **necessary** conditions for all differentiable functions (i.e. f, g_i, h_j) with optimality points $\hat{\mathbf{w}}$ and $(\hat{\lambda}, \hat{\mu})$ are:
 - $\nabla f(\hat{\mathbf{w}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{w}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{w}}) = 0$
 - $g_i(\hat{\mathbf{w}}) \leq 0; 1 \leq i \leq m$
 - $\hat{\lambda}_i \geq 0; 1 \leq i \leq m$
 - $\hat{\lambda}_i g_i(\hat{\mathbf{w}}) = 0; 1 \leq i \leq m$
 - $h_j(\hat{\mathbf{w}}) = 0; 1 \leq j \leq p$
- When f and $g_i, \forall i \in [1, m]$ are convex and $h_j, \forall j \in [1, p]$ are affine, KKT conditions are also **sufficient** for optimality at $\hat{\mathbf{w}}$ and $(\hat{\lambda}, \hat{\mu})$

KKT conditions for the Constrained (**Convex**) Problem

Recap Application 1: Equivalence of two forms of Ridge Regression

Equivalent Forms of Ridge Regression

- Consider the formulation in which we limit the weights of the coefficients by putting a constraint on size of the L2 norm of the weight vector:

$$\operatorname{argmin}_{\mathbf{w}} (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y})$$

$$\|\mathbf{w}\|_2^2 \leq \xi$$

- The objective function, namely $f(\mathbf{w}) = (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y})$ is strictly convex. The constraint function, $g(\mathbf{w}) = \|\mathbf{w}\|_2^2 - \xi$, is also convex.
- For convex $g(\mathbf{w})$, the set $\{\mathbf{w} | g(\mathbf{w}) \leq 0\}$, is also convex. (Why?)

Equivalent Forms of Ridge Regression

- To minimize the error function subject to constraint $\|\mathbf{w}\| \leq \xi$, we apply KKT conditions at the point of optimality \mathbf{w}^*

$$\nabla_{\mathbf{w}^*}(f(\mathbf{w}) + \lambda g(\mathbf{w})) = \mathbf{0}$$

(the first KKT condition). Here, $f(\mathbf{w}) = (\Phi\mathbf{w} - \mathbf{y})^T(\Phi\mathbf{w} - \mathbf{y})$ and, $g(\mathbf{w}) = \|\mathbf{w}\|^2 - \xi$.

- Solving we get,

$$\mathbf{w}^* = (\Phi^T\Phi + \lambda I)^{-1}\Phi^T\mathbf{y}$$

From the second KKT condition we get,

$$\|\mathbf{w}^*\|^2 \leq \xi$$

From the third KKT condition,

$$\lambda \geq 0$$

From the fourth condition

$$\lambda\|\mathbf{w}^*\|^2 = \lambda\xi$$

Equivalent Forms of Ridge Regression

- Values of \mathbf{w} and λ that satisfy all these equations would yield an optimal solution. That is, if

$$\|\mathbf{w}^*\| = \|(\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}\| \leq \xi$$

then $\lambda = 0$ is the solution. Else, for some sufficiently large value, λ will be the solution to

$$\|\mathbf{w}^*\| = \|(\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y}\| = \xi$$

Bound on λ in the regularized least square solution

- Consider,

$$(\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y} = \mathbf{w}^*$$

We multiply $(\Phi^T \Phi + \lambda I)$ on both sides and obtain,

$$\|(\Phi^T \Phi) \mathbf{w}^* + (\lambda I) \mathbf{w}^*\| = \|\Phi^T \mathbf{y}\|$$

Using the triangle inequality we obtain,

$$\|(\Phi^T \Phi) \mathbf{w}^*\| + (\lambda) \|\mathbf{w}^*\| \geq \|(\Phi^T \Phi) \mathbf{w}^* + (\lambda I) \mathbf{w}^*\| = \|\Phi^T \mathbf{y}\|$$

- By the Cauchy Schwarz inequality, $\|(\Phi^T \Phi) \mathbf{w}^*\| \leq \alpha \|\mathbf{w}^*\|$ for some $\alpha = \|(\Phi^T \Phi)\|$. Substituting in the previous equation,

$$(\alpha + \lambda) \|\mathbf{w}^*\| \geq \|\Phi^T \mathbf{y}\|$$

i.e.

$$\lambda \geq \frac{\|\Phi^T \mathbf{y}\|}{\|\mathbf{w}^*\|} - \alpha$$

Note that when $\|\mathbf{w}^*\| \rightarrow \mathbf{0}$, $\lambda \rightarrow \infty$. (Any intuition?) Using $\|\mathbf{w}^*\|^2 \leq \xi$ we get,

Bound on λ in the regularized least square solution

$\|(\Phi^T \Phi) \mathbf{w}^*\| \leq \alpha \|\mathbf{w}^*\|$ for some α for finite $\|(\Phi^T \Phi) \mathbf{w}^*\|$. Substituting in the previous equation,

$$(\alpha + \lambda) \|\mathbf{w}^*\| \geq \|\Phi^T \mathbf{y}\|$$

i.e.

$$\lambda \geq \frac{\|\Phi^T \mathbf{y}\|}{\|\mathbf{w}^*\|} - \alpha$$

Note that when $\|\mathbf{w}^*\| \rightarrow 0, \lambda \rightarrow \infty$. (Any intuition?) Using $\|\mathbf{w}^*\|^2 \leq \xi$ we get,

$$\lambda \geq \frac{\|\Phi^T \mathbf{y}\|}{\sqrt{\xi}} - \alpha$$

This is not the exact solution of λ but the bound proves the existence of λ for some ξ and Φ .

The Resultant alternative objective function

Substituting $g(\mathbf{w}) = \|\mathbf{w}\|^2 - \xi$, in the first KKT equation considered earlier:

$$\nabla_{\mathbf{w}^*}(f(\mathbf{w}) + \lambda \cdot (\|\mathbf{w}\|^2 - \xi)) = \mathbf{0}$$

This is equivalent to solving

$$\min(\|\Phi\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2)$$

for the same choice of λ . This form of **regularized** ridge regression is the **penalized ridge regression**.

KKT conditions for the Constrained (**Convex**) Problem

Application 2: SVR and its Dual

KKT and Dual for SVR

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

- Let's consider the lagrange multipliers α_i , α_i^* , μ_i and μ_i^* corresponding to the above-mentioned constraints.
- The Lagrange Function is

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

- Let's consider the lagrange multipliers α_i , α_i^* , μ_i and μ_i^* corresponding to the above-mentioned constraints.

- The Lagrange Function is $L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$
$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i \left(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i \right) +$$
$$\sum_{i=1}^m \alpha_i^* \left(b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^* \right) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t. \mathbf{w} ,

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t. \mathbf{w} ,
$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0 \text{ i.e., } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$
- Differentiating the Lagrangian w.r.t. ξ_i ,

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t. \mathbf{w} ,
 $\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0$ i.e., $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$
- Differentiating the Lagrangian w.r.t. ξ_i ,
 $C - \alpha_i - \mu_i = 0$ i.e., $\alpha_i + \mu_i = C$
- Differentiating the Lagrangian w.r.t ξ_i^* ,

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t. \mathbf{w} ,
 $\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0$ i.e., $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$
- Differentiating the Lagrangian w.r.t. ξ_i ,
 $C - \alpha_i - \mu_i = 0$ i.e., $\alpha_i + \mu_i = C$
- Differentiating the Lagrangian w.r.t ξ_i^* ,
 $\alpha_i^* + \mu_i^* = C$
- Differentiating the Lagrangian w.r.t b ,

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t. \mathbf{w} ,
 $\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0$ i.e., $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$
- Differentiating the Lagrangian w.r.t. ξ_i ,
 $C - \alpha_i - \mu_i = 0$ i.e., $\alpha_i + \mu_i = C$
- Differentiating the Lagrangian w.r.t ξ_i^* ,
 $\alpha_i^* + \mu_i^* = C$
- Differentiating the Lagrangian w.r.t b ,
 $\sum_i (\alpha_i^* - \alpha_i) = 0$
- Complimentary slackness:

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

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

Conclusions from the KKT conditions:

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

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

KKT conditions

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

Conclusions from the KKT conditions:

$$\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0$$

and

$$\alpha_i^*(b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) = 0$$

$\Rightarrow ?$

Conclusions from the KKT conditions:

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

$$(C - \alpha_i)\xi_i = 0 \Rightarrow ?$$

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

$$(C - \alpha_i^*)\xi_i^* = 0 \Rightarrow ?$$

For Support Vector Regression, since the original objective and the constraints are convex, any $(\mathbf{w}, b, \alpha, \alpha^*, \mu, \mu^*, \xi, \xi^*)$ that satisfy the necessary KKT conditions gives optimality (conditions are also sufficient)

Some observations

- $\alpha_i, \alpha_i^* \geq 0, \mu_i, \mu_i^* \geq 0, \alpha_i + \mu_i = C$ and $\alpha_i^* + \mu_i^* = C$
Thus, $\alpha_i, \mu_i, \alpha_i^*, \mu_i^* \in [0, C], \forall i$
- If $0 < \alpha_i < C$, then $0 < \mu_i < C$
(as $\alpha_i + \mu_i = C$)
- $\mu_i \xi_i = 0$ and $\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0$ are complementary slackness conditions

So $0 < \alpha_i < C \Rightarrow \xi_i = 0$ and $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b = \epsilon + \xi_i = \epsilon$

- All such points lie on the boundary of the ϵ band
- Using any point \mathbf{x}_j (that is with $\alpha_j \in (0, C)$) on margin, we can recover b as:
$$b = y_j - \mathbf{w}^\top \phi(\mathbf{x}_j) - \epsilon$$

KKT Conditions, Duality, SVR Dual

KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

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

For Support Vector Regression, since the original objective and the constraints are convex, any $(\mathbf{w}, b, \alpha, \alpha^*, \mu, \mu^*, \xi, \xi^*)$ that satisfy the necessary KKT conditions gives optimality (conditions are also sufficient)

Some observations

- $\alpha_i, \alpha_i^* \geq 0, \mu_i, \mu_i^* \geq 0, \alpha_i + \mu_i = C$ and $\alpha_i^* + \mu_i^* = C$
Thus, $\alpha_i, \mu_i, \alpha_i^*, \mu_i^* \in [0, C], \forall i$
- If $0 < \alpha_i < C$, then $0 < \mu_i < C$
(as $\alpha_i + \mu_i = C$)
- $\mu_i \xi_i = 0$ and $\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0$ are complementary slackness conditions

So $0 < \alpha_i < C \Rightarrow \xi_i = 0$ and $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b = \epsilon + \xi_i = \epsilon$

- All such points lie on the boundary of the ϵ band
- Using any point \mathbf{x}_j (that is with $\alpha_j \in (0, C)$) on margin, we can recover b as:
$$b = y_j - \mathbf{w}^\top \phi(\mathbf{x}_j) - \epsilon$$

Support Vector Regression

Dual Objective

Weak Duality

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$
- By weak duality theorem, we have:
$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t. $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$, and
 $\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$, and
 $\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$
- The above is true for any $\alpha_i, \alpha_i^* \geq 0$ and $\mu_i, \mu_i^* \geq 0$
- Thus,

Weak Duality

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$

- By weak duality theorem, we have:

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq L^*(\alpha, \alpha^*, \mu, \mu^*)$$

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

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

$$\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$$

- The above is true for any $\alpha_i, \alpha_i^* \geq 0$ and $\mu_i, \mu_i^* \geq 0$

- Thus,

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

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

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

$$\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$$

Dual objective

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$
- Assume: In case of SVR, we have a strictly convex objective and linear constraints
 \Rightarrow KKT conditions are necessary and sufficient and strong duality holds (for $\alpha, \alpha^* \geq 0$):

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) = \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t. $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$, and
 $\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$, and
 $\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$

- This value is precisely obtained at the $\{\hat{\mathbf{w}}, \hat{b}, \hat{\xi}, \hat{\xi}^*, \hat{\alpha}, \hat{\alpha}^*, \hat{\mu}, \hat{\mu}^*\}$ that satisfies the necessary (and sufficient) KKT optimality conditions [**KKT Constraint Set**]
- Given strong duality, we can equivalently solve: $\max_{\hat{\alpha}, \hat{\alpha}^*, \hat{\mu}, \hat{\mu}^*} L^*(\hat{\alpha}, \hat{\alpha}^*, \hat{\mu}, \hat{\mu}^*)$