Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 10 - Optimization Foundations Applied to Regression
Formulations

1. Is there a probabilistic interpretation?
   - Gaussian Error, Maximum Likelihood Estimate
2. Addressing overfitting
   - Bayesian and Maximum Aposteriori 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

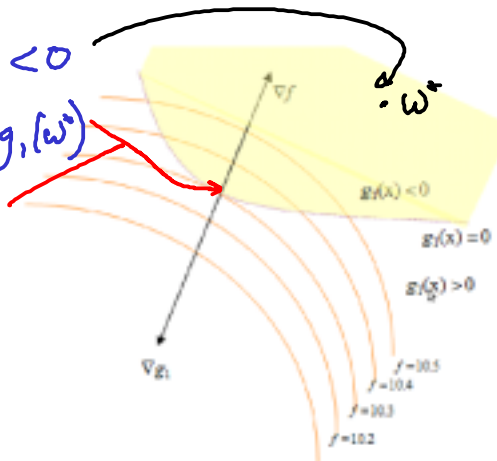**Q.** *How to solve such constrained problems?*

**A.** Canonical example:

*Minimize $f(\mathbf{w})$ s.t. $g_1(\mathbf{w}) \leq 0$* $\qquad\qquad$ (1)

At $w^*$ (optimal pt)

$\nabla f(w^*)$ & $g_1(w^*) < 0$

$\underline{\underline{OR}}$ $\quad \nabla f(w^*) = -\lambda \nabla g_1(w^*)$
$\qquad\qquad g_1(w^*) = 0$

- If $\mathbf{w}^*$ is on the boundary of $g_1$, *i.e.*, if $g_1(\mathbf{w}^*) = 0$,

$$\nabla f(\mathbf{w}^*) = -\lambda \nabla g_1(\mathbf{w}^*) \text{ for some } \lambda \geq 0$$

- **Intuition:** If the above didn't hold, then we would have $\nabla f(\mathbf{w}^*) = \lambda_1 \nabla g_1(\mathbf{w}^*) + \lambda_2 \nabla_\perp g_1(\mathbf{w}^*)$, where, by moving in direction[1] $\pm \nabla_\perp g_1(\mathbf{w}^*)$ ( or $-\nabla g_1(\mathbf{w}^*)$), we remain on boundary $g_1(\mathbf{w}^*) = 0$, ( or within $g_1(\mathbf{w}^*) \leq 0$) while decreasing the value of $f$, which is not possible at the point of optimality.

- Thus, at the point of optimality[2], for some $\lambda \geq 0$,

$$\text{Either } g_1(\mathbf{w}^*) < 0 \quad \& \quad \nabla f(\mathbf{w}^*) = 0 \quad (\lambda = 0)$$
$$\text{Or } g_1(\mathbf{w}^*) = 0 \quad \& \quad \nabla f(\mathbf{w}^*) = -\lambda \nabla g_1(\mathbf{w}^*)$$

*(handwritten annotations:)*

Summary:

$L(\omega) = f(\omega) + \lambda g_1(\omega)$

$\nabla L(\omega^*) = 0$ (2)

$\left( \begin{array}{l} \lambda g_1(\omega^*) = 0 \\ g_1(\omega^*) \leq 0 \end{array} \right)$ (3)

Comp. Slackness

[1] $\nabla_\perp g_1(\mathbf{w}^*)$ is the direction orthogonal to $\nabla g_1(\mathbf{w}^*)$
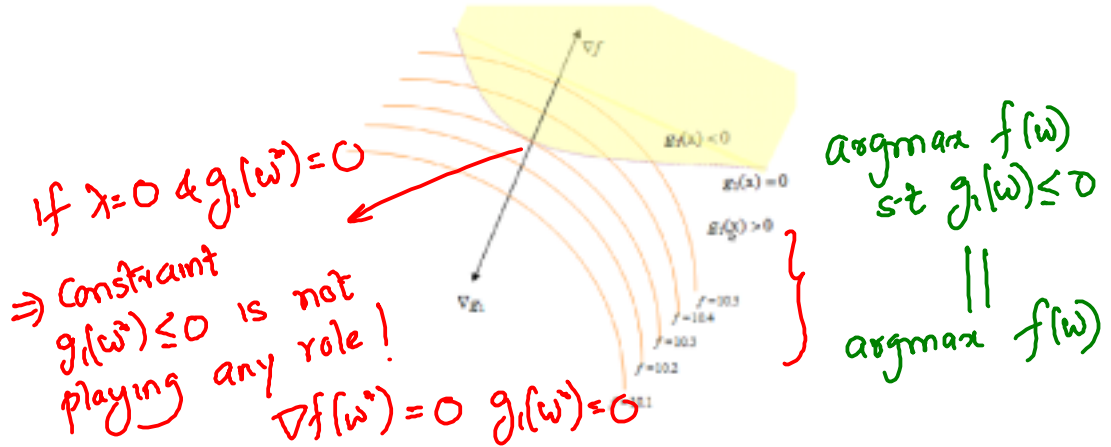[2] Section 4.4, pg-72: cs725/notes/BasicsOfConvexOptimization.pdf

Figure 2: Two conditions under which a minimum can occur: a) When the minimum is on the constraint function boundary, in which case the gradients are in opposite directions; b) When point of minimum is inside the constraint space (shown in yellow shade), in which case $\nabla f(\mathbf{w}^*) = \mathbf{0}$.

- The first condition occurs when minima lies on the boundary of function $g$. In this case, gradient vectors corresponding to the functions $f$ and $g$, at $\mathbf{w}^*$, point in opposite directions barring multiplication by a real constant.
- Second condition represents the case that point of minimum lies inside the constraint space. This space is shown shaded in Figure 1. Clearly, for this case, $\nabla f(\mathbf{w}) = \mathbf{0}$.
- An Alternative Representation: $\nabla L(\mathbf{w}, \lambda) = 0$ for some $\lambda \geq 0$ where

$$L(\mathbf{w}, \lambda) = f(\mathbf{w}) + \lambda \mathbf{g}(\mathbf{w}); \lambda \in \mathbb{R}$$

*(handwritten annotations: "original objective" pointing to $f(\mathbf{w})$; "penalized constraint" pointing to $\lambda \mathbf{g}(\mathbf{w})$)*

is called the lagrange function which has objective function augmented by weighted sum of constraint functions

For a convex objective and constraint function, the minima, $\mathbf{w}^*$, can satisfy one of the following two conditions:

1. $g(\mathbf{w}^*) = \mathbf{0}$ and $\nabla f(\mathbf{w}^*) = -\lambda \nabla \mathbf{g}(\mathbf{w}^*)$
2. $g(\mathbf{w}^*) < \mathbf{0}$ and $\nabla f(\mathbf{w}^*) = \mathbf{0}$

- Here, we wish to penalize higher magnitude coefficients, hence, we wish $g(\mathbf{w})$ to be negative while minimizing the lagrangian. In order to maintain such direction, we must have $\lambda \geq 0$. Also, for solution $\mathbf{w}$ to be feasible, $\nabla g(\mathbf{w}) \leq \mathbf{0}$.

- Due to complementary slackness condition, we further have $\lambda g(\mathbf{w}) = \mathbf{0}$, which roughly suggests that the lagrange multiplier is zero unless constraint is active at the minimum point. As $\mathbf{w}$ minimizes the lagrangian $L(\mathbf{w}, \lambda)$, gradient must vanish at this point and hence we have $\nabla f(\mathbf{w}) + \lambda \nabla \mathbf{g}(\mathbf{w}) = \mathbf{0}$

KKT Conditions, Duality, SVR Dual

- The general optimization problem we consider with (convex) inequality and (linear) equality constraints is:

$$\min_{\mathbf{w}} f(\mathbf{w})$$

*Illustrate through ← Quadratic fn (SVR)*

$$\|A\omega - b\|_2^2$$

- The general optimization problem we consider with (convex) inequality and (linear) equality constraints is:

$$\min_{\mathbf{w}} f(\mathbf{w})$$

Set of inequalities
eg: linear $B\mathbf{w} \geq C$ } subject to $g_i(\mathbf{w}) \leq 0; 1 \leq i \leq m$

Set of equalities
eg: linear $D\mathbf{w} = f$ } $h_j(\mathbf{w}) = 0; 1 \leq j \leq p$

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

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

Goal:
Minimize $L$

Diff penalty with each equality

Intuition (Necessary at opt)

① Since we need
$g_i(\omega) \downarrow \leq 0, \ \lambda_i \geq 0$

Penalty associated with each inequality
$g_i(\omega) \leq 0$

$\nabla f(\omega^*) = -\sum_i \lambda_i \nabla g_i(\omega^*) - \sum_j \mu_j \nabla h_j(\omega^*)$

② Since we need
$h_j(\omega) \approx 0$, sign of $\mu_j$ does not matter

$g_i(\omega^*) \leq 0$

$h_j(\omega^*) = 0$

$\lambda_i g_i(\omega^*) = 0$

Karush Kuhn Tucker

- 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:

① $\nabla f(\hat{\mathbf{w}}) + \sum_{i=1}^{m} \hat{\lambda}_i \nabla g_i(\hat{\mathbf{w}}) + \sum_{j=1}^{p} \hat{\mu}_j \nabla h_j(\hat{\mathbf{w}}) = 0$ → Gradient equality

② $g_i(\hat{\mathbf{w}}) \leq 0; 1 \leq i \leq m$ → ineq (original)

③ $\hat{\lambda}_i \geq 0; 1 \leq i \leq m$ → positivity

④ $\hat{\lambda}_i g_i(\hat{\mathbf{w}}) = 0; 1 \leq i \leq m$ → comp slackness

⑤ $h_j(\hat{\mathbf{w}}) = 0; 1 \leq j \leq p$ → eq (original)

Let us apply to constrained ridge regression

$$\min_{\omega} \; \|\phi\omega - y\|_2^2$$

$$s.t \quad \|\omega\|_2^2 \leq \theta \rightarrow \lambda$$

Looks like soln to penalized ridge regression

$$\hat{\omega} = (\Phi^T\Phi + \hat{\lambda}I)^{-1}\Phi^T y$$

$$L(\omega, \lambda) = \|\phi\omega - y\|_2^2 + \lambda\left(\|\omega\|_2^2 - \theta\right)$$

① $\nabla L(\hat{\omega}, \hat{\lambda}) = 0 \Rightarrow 2\Phi^T\Phi\hat{\omega} - 2\Phi^T y + 2\lambda\hat{\omega} = 0$ $\Big\}$ $\Big\downarrow \hat{\lambda}$ s.t

② $\|\hat{\omega}\|_2^2 \leq \theta$  ③ $\hat{\lambda} \geq 0$  ④ $\hat{\lambda}\left(\|\hat{\omega}\|_2^2 - \theta\right) = 0$ $\Big\}$ $\hat{\lambda} = f(\theta)$

$\hat{\lambda}$ s.t if $\|(\phi^T\phi)^{-1}\phi^T y\|_2^2 \leq \theta$ then $\hat{\lambda} = 0$ else smallest $\hat{\lambda} > 0$

s.t $\|(\phi^T\phi + \hat{\lambda}I)^{-1}\phi^T y\|_2^2 = \theta$

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

*(handwritten annotation)* $h_j \leq 0$ } both convex, $-h_j \leq 0$ } is linear

- 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 $(\hat{\lambda}, \hat{\mu})$

- With $\mathbf{w} \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p$, Lagrangian is:

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

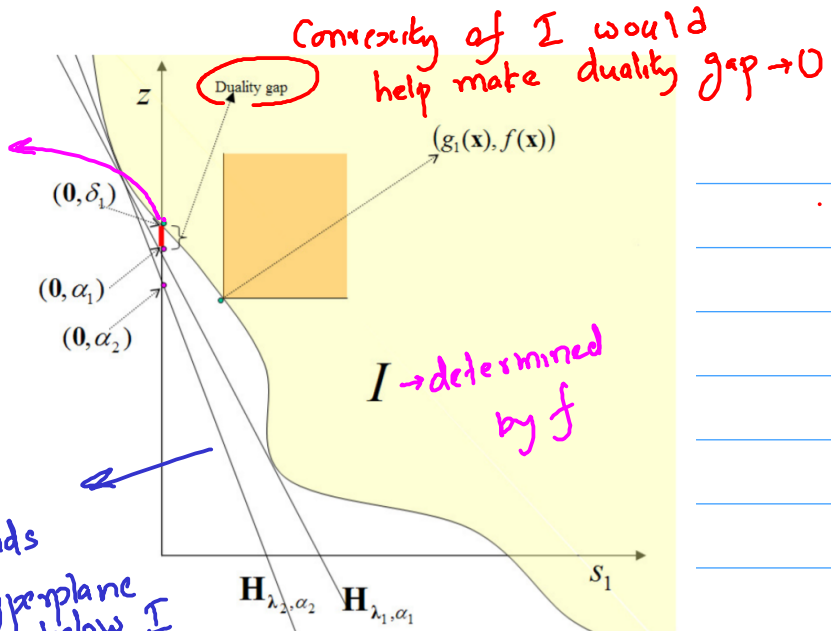- Lagrange dual function is minimum of Lagrangian over $\mathbf{w}$.

$$L^*(\lambda, \mu) = \min_{\omega} L(\omega, \lambda, \mu) \leq f(\omega) \quad \text{s.t} \quad g_i(\omega) \leq 0$$
$$h_j(\omega) = 0$$

$$\max_{\lambda, \mu} L^*(\lambda, \mu) \leq f(\omega) \quad \text{s.t} \quad g_i(\omega) \leq 0$$
$$h_j(\omega) = 0$$

Red gap in image will always lie above Hyperplane

Convexity of $I$ would help make duality gap $\to 0$

Duality gap

$z$

$(g_1(\mathbf{x}), f(\mathbf{x}))$

$(\mathbf{0}, \delta_1)$

$(\mathbf{0}, \alpha_1)$

$(\mathbf{0}, \alpha_2)$

y intercept of $I$ $\equiv$ optimal value of $f$ subject to constraints

$I \to$ determined by $f$

$L^*(\lambda, \mu)$ corresponds to a Hyperplane below $I$

$\mathbf{H}_{\lambda_2, \alpha_2}$    $\mathbf{H}_{\lambda_1, \alpha_1}$

$s_1$

- With $\mathbf{w} \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p$, Lagrangian is:

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

- Lagrange dual function is minimum of Lagrangian over $\mathbf{w}$.

$$L^*(\lambda, \mu) = \min_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu)$$

- The Dual Optimization Problem is to maximize Lagrange dual function $L^*(\lambda, \mu)$ over $(\lambda, \mu)$

$$\hat{\lambda}, \hat{\mu} = \underset{\lambda, \mu}{\arg\max} \; L^*(\lambda, \mu) \equiv \text{Push the Hyperplane as upward as possible}$$

# Lagrangian Duality and KKT conditions

- With $\mathbf{w} \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p$, Lagrangian is:

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

- Lagrange dual function is minimum of Lagrangian over $\mathbf{w}$.

$$L^*(\lambda, \mu) = \min_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu)$$

- The Dual Optimization Problem is to maximize Lagrange dual function $L^*(\lambda, \mu)$ over $(\lambda, \mu)$

$$\underset{\lambda,\mu}{\operatorname{argmax}} \ L^*(\lambda, \mu) = \underset{\lambda,\mu}{\operatorname{argmax}} \ \min_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu)$$

*tends to = for convex prob.*

$$\max_{\lambda \mu} L^*(\lambda, \mu) = \max_{\lambda \mu} \min_{\omega} L(\omega, \lambda, \mu) \leq \min_{\omega} f(\omega)$$
$$g_i(\omega) \leq 0$$
$$h_i(\omega) \leq 0$$

- The dual function yields lower bound for minimizer of the primal formulation.
- Max of dual function $L^*(\lambda, \mu)$ over $(\lambda, \mu)$ is also therefore a lower bound

$$L^*(\lambda, \mu) \leq f(\omega) \quad \forall \; \omega, \lambda, \mu$$

$$\lambda \geq 0 \qquad st \quad g_i(\omega) \leq 0$$

$$h_j(\omega) = 0$$

max over $\lambda, \mu$

min over $\omega$

- The dual function yields lower bound for minimizer of the primal formulation. ✓
- Max of dual function $L^*(\lambda, \mu)$ over $(\lambda, \mu)$ is also therefore a lower bound

$$\max_{\lambda,\mu} L^*(\lambda, \mu) = \max_{\lambda,\mu} \min_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu) \leq L(\mathbf{w}, \lambda, \mu)$$

- **Duality Gap:** The gap between primal and dual solutions. In the KKT conditions, $\hat{\mathbf{w}}$ correspond to primal optimal and $(\hat{\lambda}, \hat{\mu})$ to dual optimal points $\Rightarrow$ Duality gap is $f(\hat{\mathbf{w}}) - L^*(\hat{\lambda}, \hat{\mu})$
- Duality gap characterizes suboptimality of the solution and can be approximated by $f(\mathbf{w}) - L^*(\lambda, \mu)$ for any feasible $\mathbf{w}$ and corresponding $\lambda$ and $\mu$

if $\hat{\lambda}, \hat{\mu}, \hat{\omega}$ is soln to KKT condition then
$f(\hat{\omega}) - L^*(\hat{\lambda}, \hat{\mu})$ is gap.
$= 0$ under convexity

Most imp

- The dual function yields lower bound for minimizer of the primal formulation.
- Max of dual function $L^*(\lambda, \mu)$ over $(\lambda, \mu)$ is also therefore a lower bound

$$\max_{\lambda,\mu} L^*(\lambda, \mu) = \max_{\lambda,\mu} \min_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu) \leq L(\mathbf{w}, \lambda, \mu)$$

- **Duality Gap:** The gap between primal and dual solutions. In the KKT conditions, $\hat{\mathbf{w}}$ correspond to primal optimal and $(\hat{\lambda}, \hat{\mu})$ to dual optimal points $\Rightarrow$ Duality gap is $f(\hat{\mathbf{w}}) - L^*(\hat{\lambda}, \hat{\mu})$
- Duality gap characterizes suboptimality of the solution and can be approximated by $f(\mathbf{w}) - L^*(\lambda, \mu)$ for any feasible $\mathbf{w}$ and corresponding $\lambda$ and $\mu$
- When functions $f$ and $g_i, \forall i \in [1, m]$ are convex and $h_j, \forall j \in [1, p]$ are affine, Karush-Kuhn-Tucker (KKT) conditions are both necessary and sufficient for points to be both primal and dual optimal with zero duality gap.

Elaboration on equivalence of penalized & constrained forms of ridge regression (continued from page 14)

- 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:

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

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

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

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

(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

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

- Values of **w** and $\lambda$ 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, $\lambda$ will be the solution to

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

- 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,

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

  Using the triangle inequality we obtain,

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

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

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

  i.e.

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

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

$\|(\Phi^T\Phi)\mathbf{w}^*\| \leq \alpha\|\mathbf{w}^*\|$ for some $\alpha$ 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 $\|\mathbf{w}^*\| \to 0, \lambda \to \infty$. (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 $\lambda$ but the bound proves the existence of $\lambda$ for some $\xi$ and $\Phi$.

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 $\lambda$. This form of **regularized** ridge regression is the **penalized ridge regression**.