

Lecture 07- Convex and Constrained Optimization and Regression

Instructor: Prof. Ganesh Ramakrishnan

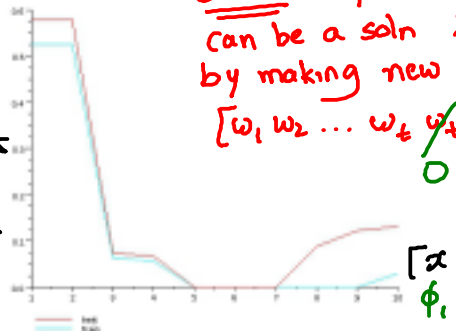
Agenda

- Overfitting
- Regularized Regression and Constrained Convex Optimization
- Support Vector Regression

Redundant ϕ and Overfitting

(Similar to standard dev)
 Regression error
 Root mean sq. error

$$\sqrt{\frac{\sum (y_i - \hat{\phi}(x_i))^2}{n}}$$



Claim: Any soln for t
 can be a soln for $t' > t$
 by making new coeffs = 0

$$[\omega_1 \omega_2 \dots \omega_t \quad \cancel{\omega_{t+1}} \dots \cancel{\omega_{t'}}]$$

$$[x \quad x^2 \quad \dots \quad x^t]$$

$$[\phi_1 \phi_2 \quad \dots \quad \phi_t]$$

Figure: train-RMS and test-RMS values vs t (degree of polynomial) graph

- Too many bends ($t=9$ onwards) in curve \equiv high values of some w_i 's
- Train and test errors differ significantly

Constrained Least Squares Linear Regression

Find

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|^2 \quad \text{s.t.} \quad \underbrace{\|\mathbf{w}\|_p^p}_{\Omega(\mathbf{w})} \leq \zeta, \quad (1)$$

where

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

Why? Suppose $p \rightarrow 0$, Note: $x^0 = 1$ except if $x = 0$

$\lim_{p \rightarrow 0} \|\mathbf{w}\|_p = \# \text{ of non-zero } w_i\text{'s}$

$0^p = 0$

eg of feature selection

Also called support/cardinality of \mathbf{w}

For other p 's let us look at level curves of $\|\mathbf{w}\|_p$

p-Norm level curves

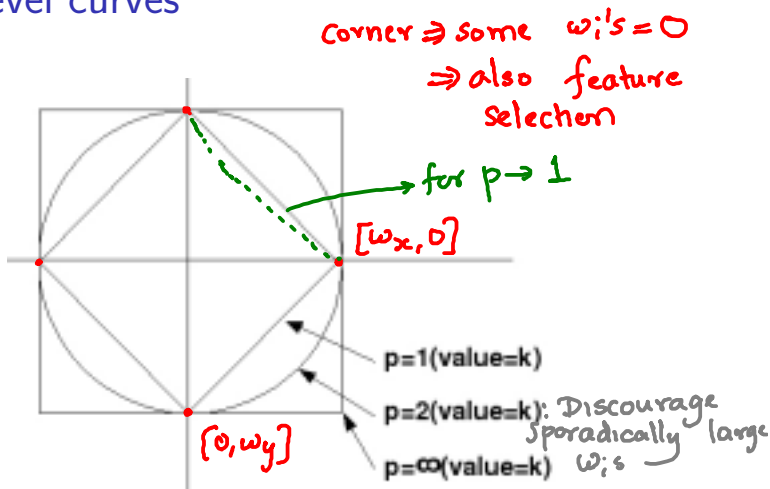


Figure: p-Norm curves for constant norm value and different p

Convex Optimization Problem

- Formally, a convex optimization problem is an optimization problem of the form

$$\text{minimize } f(\mathbf{x}) \Leftrightarrow \|\phi\omega - \mathbf{y}\|_2^2 \quad (3)$$

$$\text{subject to } c \in C \Leftrightarrow \|\omega\|_p^p \leq \xi \quad (4)$$

where f is a convex function, C is a convex set, and \mathbf{x} is the optimization variable.

- An improved form of the above would be

$$\text{minimize } f(\mathbf{x}) \quad (5)$$

$$g_i(\mathbf{x}) = \|\omega\|_p^p - \xi \leq 0, \quad i = 1, \dots, m \quad (6)$$

$$h_i(\mathbf{x}) = 0, \quad i = 1, \dots, p \quad (7)$$

where f is a convex function, g_i are convex functions, and h_i are affine functions, and \mathbf{x} is the vector of optimization variables.

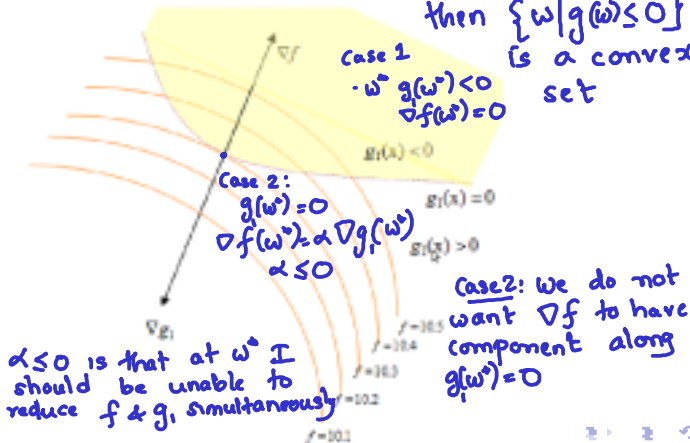
Constrained convex problems

Q. How to solve constrained problems of the above-mentioned type?

A. General problem format :

$$\text{Minimize } f(\mathbf{w}) \text{ s.t. } g(\mathbf{w}) \leq 0 = f(\mathbf{w}^*) \quad (8)$$

If $g(\mathbf{w})$ is a convex fn then $\{\mathbf{w} | g(\mathbf{w}) \leq 0\}$ is a convex set



Constrained Convex Problems

- At the point of optimality,

$$\text{Either } g(\mathbf{w}^*) < 0 \quad \& \quad \nabla f(\mathbf{w}^*) = 0 \quad (9)$$

$$\text{Or } g(\mathbf{w}^*) = 0 \quad \& \quad \nabla f(\mathbf{w}^*) = \alpha \nabla g(\mathbf{w}^*) \quad (10)$$

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

α is a Lagrange multiplier

$$\nabla f(\mathbf{w}^*) = \alpha \nabla g(\mathbf{w}^*) \quad (11)$$

(Duality Theory) (Cite : Section 4.4, pg-72)¹

- **Intuition:** If the above didn't hold, then we would have $\nabla f(\mathbf{w}^*) = \alpha_1 \nabla g(\mathbf{w}^*) + \alpha_2 \nabla_{\perp} g(\mathbf{w}^*)$, where by moving in direction $\pm \nabla_{\perp} g(\mathbf{w}^*)$, we remain on boundary $g(\mathbf{w}^*) = 0$, while decreasing/increasing value of f , which is not possible at the point of optimality.

"Regularized" Linear Regression

- We limit the weights of the coefficients by putting a constraint on size of the L2 norm of the weight vector

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

$$g_1(\mathbf{w}) = \|\mathbf{w}\|_2^2 - \xi \leq 0 \Leftrightarrow \|\mathbf{w}\|_2^2 \leq \xi$$

- 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?)

Duality and KKT conditions

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

- 1 $g(w^*) = 0$ and $\nabla f(w^*) = \alpha \nabla g(w^*)$
- 2 $g(w^*) < 0$ and $\nabla f(w^*) = 0$

Require α s.t.

$$\|(\Phi^T \Phi - 2\alpha I)^{-1} \Phi^T y\|_2^2 = \xi$$

$$\stackrel{\text{I}}{\cong} \|w^*\|_2^2 = \xi \quad \& \quad (\Phi^T \Phi) w^* - \Phi^T y = 2\alpha w^* \Rightarrow w^* = (\Phi^T \Phi - 2\alpha I)^{-1} \Phi^T y$$

$$\stackrel{\text{OR}}{\cong} \|w^*\|_2^2 < \xi \quad \& \quad (\Phi^T \Phi) w^* - \Phi^T y = 0 \Rightarrow w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

Requires: $\|(\Phi^T \Phi)^{-1} \Phi^T y\|_2^2 < \xi$

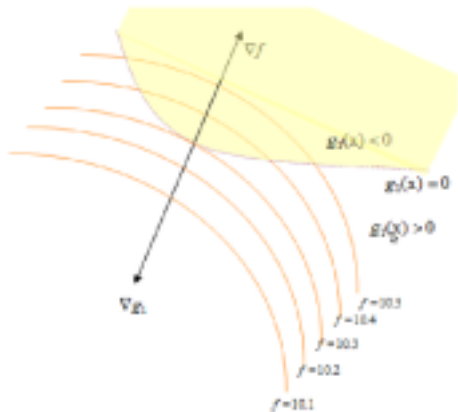


Figure: Two conditions when a minima can occur: a) When the minima is on the constraint function boundary, in which case the gradients are along the same direction ;b) When minima is inside the constraint space (shown in yellow shade), in which case $\nabla f(\mathbf{w}^*) = \mathbf{0}$.

Duality and KKT conditions

- This fact can be easily visualized from the previous figure. As we can see, the first condition occurs when minima lies on the boundary of function g . In this case, gradient vectors corresponding to the function f and the function g , at \mathbf{w}^* , point in the same direction barring multiplication by a real constant.
- Second condition depicts the case when minima lies inside the constraint space. This space is shown shaded in Figure 1. Clearly, for this case $\nabla f(\mathbf{w}^*) = \mathbf{0}$ for minima to occur. This primal problem can be converted to dual using the lagrange multiplier. According to which, we can convert this problem to the objective function augmented by weighted sum of constraint functions in order to get the corresponding lagrangian.

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

Duality and KKT conditions

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

Duality and KKT conditions

- In general, optimization problem with inequality and equality constraints might be depicted in the following manner:

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

$$\text{subject to } \underline{g_i(\mathbf{w})} \leq \underline{0}; 1 \leq i \leq m$$

inequality

$$\underline{h_j(\mathbf{w})} = \underline{0}; 1 \leq j \leq p$$

equality

Duality and KKT conditions

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

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

Handwritten annotations: $\lambda_i \geq 0$ (with an arrow pointing to λ_i), and Multiplicators (with arrows pointing to λ_i and μ_j). The term $L(\mathbf{w}, \lambda, \mu)$ is circled in red.

- Lagrange dual function is the minimum value of the lagrangian over $\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p$.

$$L^*(\lambda, \mu) = \operatorname{argmin}_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu)$$

Necessary condition for $\min_{\mathbf{w}} L(\mathbf{w}, \lambda, \mu)$ at $\hat{\mathbf{w}}$: $\nabla L(\hat{\mathbf{w}}) = 0 \Rightarrow \nabla f(\hat{\mathbf{w}}) + \sum_i \lambda_i \nabla g_i(\hat{\mathbf{w}}) + \sum_j \mu_j \nabla h_j(\hat{\mathbf{w}}) = 0$

Duality and KKT conditions

- The **dual function** yields lower bound for minimizer of the **primal** formulation.
 $L^*(\lambda, \mu) = \min_{\omega} L(\omega, \lambda, \mu)$
 $L^*(\lambda, \mu) \leq \min_{\omega} f(\omega) \text{ s.t. } g_i(\omega) \leq 0, h_j(\omega) = 0$
- Max of dual function $L^*(\lambda, \mu)$ over (λ, μ) is also therefore a lower bound
 $\max_{\lambda \geq 0} L^*(\lambda, \mu) \leq \min_{\omega} f(\omega) \text{ s.t. } g_i(\omega) \leq 0, h_j(\omega) = 0$
- The gap between primal and dual solutions is the duality gap,
- Duality gap characterizes suboptimality of the solution.

$$f(\mathbf{w}) - L^*(\lambda, \mu)$$

- 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.

Duality and KKT conditions

For above 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 may be given in the following manner:

$$\nabla_{\hat{\mathbf{w}}} \left(f(\hat{\mathbf{w}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{w}}) + \sum_{j=1}^p \hat{\mu}_j h_j(\hat{\mathbf{w}}) \right) = 0 \quad \checkmark$$

- $g_i(\hat{\mathbf{w}}) \leq 0; 1 \leq i \leq m \quad \checkmark$

- $\hat{\lambda}_i \geq 0; 1 \leq i \leq m \quad \checkmark$

- $\hat{\lambda}_i g_i(\hat{\mathbf{w}}) = 0; 1 \leq i \leq m$

- $h_j(\hat{\mathbf{w}}) = 0; 1 \leq j \leq p \quad \checkmark$

→ Complementary slackness condition

If $\hat{\mathbf{w}}$ on bndry $\Rightarrow g_i(\hat{\mathbf{w}}) = 0$

If $\hat{\mathbf{w}}$ NOT on bndry $\Rightarrow \nabla f(\hat{\mathbf{w}}) = 0$

$$\stackrel{!}{=} \nabla f(\hat{\mathbf{w}}) = 0 \cdot \nabla g_i(\hat{\mathbf{w}}) \\ \hat{\lambda}_i = 0$$