

Constrained least squares as the limit solution of a regularized least squares problem

Saurabh V. Pendse^{*1}

¹ Department of Computer Science and Engineering
The Maharaja Sayajirao University of Baroda, Gujarat, India

May 23, 2010

Abstract

One of the most flexible techniques for function approximation in statistics is the basis function based linear regression. By choosing a sufficiently flexible set of basis functions, one can model complicated functions accurately. Smoothness or complexity of the resulting functions can be controlled by introducing a regularization term in the objective function. In this technical note, we describe how one can solve the function estimation problem in the presence of simple linear constraints. In particular, we show that the constrained least squares solution can be obtained as the limiting solution of a particular regularized least squares problem.

1 Notation

Vectors will be denoted using lower case letters in bold face, for example \mathbf{x}_i, \mathbf{y} . Matrices will be denoted using upper case letters in bold face, for example

^{*}To whom correspondence should be addressed. e-mail: sau2pen@gmail.com

\mathbf{X} , \mathbf{A} . The j th element of a vector \mathbf{x}_i will be denoted by x_{ij} whereas the j th element of a vector \mathbf{y} will be denoted by y_j . The ij th element of a matrix \mathbf{X} will be denoted by X_{ij} . The transpose of a matrix \mathbf{A} will be denoted by \mathbf{A}^T and its inverse will be denoted by \mathbf{A}^{-1} . We will use $\mathbf{0}$ to denote a vector or matrix of all zeros whose size should be clear from context.

2 Introduction

Suppose we are given n input vectors, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ each of size $p \times 1$. Corresponding to each \mathbf{x}_i we observe an output variable y_i to get a set of n observed outputs y_1, y_2, \dots, y_n . A linear model explaining the observed y_i in terms of \mathbf{x}_i is given by:

$$y_i = \mathbf{x}_i^T \boldsymbol{\psi} + \varepsilon_i \quad (2.1)$$

where $\boldsymbol{\psi}$ is a $p \times 1$ coefficient vector and ε_i is a scalar random variable representing the unobserved noise corruption with $E(\varepsilon_i) = 0$ and $E(\varepsilon_i^2) = \sigma^2$. We also assume that ε_i is independent across i . A linear equation such as 2.1 has limited flexibility and might not be appropriate for modeling real world data. This problem is easily solved by first transforming the $p \times 1$ input vectors \mathbf{x}_i into $m \times 1$ vectors $\boldsymbol{\phi}(\mathbf{x}_i)$ where $\boldsymbol{\phi} : \mathbf{R}^p \rightarrow \mathbf{R}^m$ is a non-linear function which takes p dimensional vectors to m dimensional vectors. Each component ϕ_i of $\boldsymbol{\phi}$ represents one basis function. Next, we construct a linear model in the transformed space:

$$y_i = \boldsymbol{\phi}(\mathbf{x}_i)^T \boldsymbol{\beta} + \varepsilon_i \quad (2.2)$$

Here $\boldsymbol{\beta}$ is a $m \times 1$ coefficient vector corresponding to the m basis functions. To simplify notation define:

$$\mathbf{y} = [y_1, y_2, \dots, y_n]^T \text{ is a } n \times 1 \text{ vector} \quad (2.3)$$

$$\boldsymbol{\Phi} = [\boldsymbol{\phi}(\mathbf{x}_1), \boldsymbol{\phi}(\mathbf{x}_2), \dots, \boldsymbol{\phi}(\mathbf{x}_n)]^T \text{ is a } n \times m \text{ matrix} \quad (2.4)$$

$$\boldsymbol{\varepsilon} = [\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n]^T \text{ is a } n \times 1 \text{ vector} \quad (2.5)$$

In this work, we will assume that $n > m$ and that $\text{rank}(\boldsymbol{\Phi}) = m$. Then we can re-write 2.2 as:

$$\mathbf{y} = \boldsymbol{\Phi} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (2.6)$$

where $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $E(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T) = \sigma^2 \mathbf{I}_n$. Given \mathbf{y} and $\boldsymbol{\Phi}$ the problem is to estimate the unknown coefficient vector $\boldsymbol{\beta}$. In the least squares solution

approach, we estimate $\boldsymbol{\beta}$ by solving the optimization problem:

$$\min_{\boldsymbol{\beta}} f(\boldsymbol{\beta}) = (\mathbf{y} - \Phi\boldsymbol{\beta})^T(\mathbf{y} - \Phi\boldsymbol{\beta}) \quad (2.7)$$

Taking the gradient of 2.7 and equating it to $\mathbf{0}$ at the least squares solution point $\boldsymbol{\beta}_{LS}$ we get:

$$\frac{\partial}{\partial \boldsymbol{\beta}} f(\boldsymbol{\beta})|_{\boldsymbol{\beta}_{LS}} = 2\Phi^T (\Phi\boldsymbol{\beta}_{LS} - \mathbf{y}) = \mathbf{0} \quad (2.8)$$

Since Φ has full column rank, upon simplification, we get:

$$\boldsymbol{\beta}_{LS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad (2.9)$$

A desirable property of $\boldsymbol{\beta}_{LS}$ is that it is an unbiased estimator of $\boldsymbol{\beta}$ in the presence of noise $\boldsymbol{\varepsilon}$. Assuming that the true $\boldsymbol{\beta}$ is given, it is easy to see from 2.9 and 2.6 that $E(\boldsymbol{\beta}_{LS}) = (\Phi^T \Phi)^{-1} \Phi^T E(\mathbf{y}) = \boldsymbol{\beta}$. Since we are given only a finite number of points n , the error function $f(\boldsymbol{\beta})$ can be made arbitrarily small by selecting a sufficiently flexible set of basis functions $\phi(\mathbf{x})$. For example, any function $\phi(\mathbf{x})^T \boldsymbol{\beta}$ that passes through the given points \mathbf{y} exactly will make the objective $f(\boldsymbol{\beta}) = 0$. How can we impose conditions on the estimated function $\phi(\mathbf{x})^T \boldsymbol{\beta}$ to avoid a non-sensical noisy solution? One way is to add a penalty term to the objective function $f(\boldsymbol{\beta})$ that will encourage smoothness in the estimated function $\phi(\mathbf{x})^T \boldsymbol{\beta}$. Let $g(\mathbf{x}) = \phi(\mathbf{x})^T \boldsymbol{\beta} = \sum_{i=1}^m \phi_i(\mathbf{x}) \beta_i$ then

$$\frac{\partial^2 g(\mathbf{x})}{\partial x_k \partial x_l} = \sum_{i=1}^m \frac{\partial^2 \phi_i(\mathbf{x})}{\partial x_k \partial x_l} \beta_i \quad (2.10)$$

Using the Cauchy-Schwarz inequality we can write:

$$\left[\frac{\partial^2 g(\mathbf{x})}{\partial x_k \partial x_l} \right]^2 \leq \sum_{i=1}^m \left[\frac{\partial^2 \phi_i(\mathbf{x})}{\partial x_k \partial x_l} \right]^2 \sum_{i=1}^m \beta_i^2 \quad (2.11)$$

For a fixed set of basis functions $\phi_i(\mathbf{x})$ we can impose smoothness by forcing $\left[\frac{\partial^2 g(\mathbf{x})}{\partial x_k \partial x_l} \right]^2$ to have small values for all k and l . It is clear from 2.11 that this can be attained by making $\sum_{i=1}^m \beta_i^2 = \boldsymbol{\beta}^T \boldsymbol{\beta}$ small. This motivates the addition of a penalty term proportional to $\boldsymbol{\beta}^T \boldsymbol{\beta}$ into $f(\boldsymbol{\beta})$ in 2.7 to get a regularized least squares problem also known as ridge regression (see [1]):

$$\min_{\boldsymbol{\beta}} f_{reg}(\boldsymbol{\beta}) = (\mathbf{y} - \Phi\boldsymbol{\beta})^T(\mathbf{y} - \Phi\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \quad (2.12)$$

Here the parameter λ controls the degree of regularization or smoothness. Higher values of λ force a smoother function and lower values of λ encourage a rougher function. In addition to imposing regularization it might be desirable for β to satisfy certain linear constraints. Suppose \mathbf{G} is a $q \times m$ matrix with $q < m$ and $\text{rank}(\mathbf{G}) = q$ and let \mathbf{c} be a $q \times 1$ vector. We would like β to satisfy the constraints:

$$\mathbf{G}\beta = \mathbf{c} \quad (2.13)$$

We propose to introduce the constraints in 2.13 as "soft" constraints into the objective function in 2.12. We do not require the exact satisfaction of these soft constraints but only encourage their satisfaction. Define the modified optimization problem which encourages satisfaction of 2.13 as:

$$\min_{\beta} f_{con}(\beta) = \underbrace{(\mathbf{y} - \Phi\beta)^T(\mathbf{y} - \Phi\beta)}_{\text{model fit}} + \underbrace{\lambda\beta^T\beta}_{\text{smoothness}} + \underbrace{\theta(\mathbf{G}\beta - \mathbf{c})^T(\mathbf{G}\beta - \mathbf{c})}_{\text{soft constraints}} \quad (2.14)$$

First we obtain a solution $\beta(\lambda, \theta)$ of the problem 2.14 and then prove that $\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = \beta^*(\lambda)$ where $\beta^*(\lambda)$ solves the "hard" constraint optimization problem:

$$\min_{\beta} f_{hcon}(\beta) = \underbrace{(\mathbf{y} - \Phi\beta)^T(\mathbf{y} - \Phi\beta)}_{\text{model fit}} + \underbrace{\lambda\beta^T\beta}_{\text{smoothness}} \quad (2.15)$$

$$\text{subject to : } \underbrace{\mathbf{G}\beta = \mathbf{c}}_{\text{hard constraint}}$$

3 The soft constrained problem

Suppose $\beta(\lambda, \theta)$ is a solution of the soft constrained problem 2.14. The necessary conditions for optimality are:

$$\frac{\partial f_{con}(\beta)}{\partial \beta} \Big|_{\beta(\lambda, \theta)} = [-2\Phi^T(\mathbf{y} - \Phi\beta) + 2\lambda\beta + 2\theta\mathbf{G}^T(\mathbf{G}\beta - \mathbf{c})] \Big|_{\beta(\lambda, \theta)} = \mathbf{0} \quad (3.1)$$

Rearranging we get:

$$[\Phi^T\Phi + \lambda\mathbf{I}_m + \theta\mathbf{G}^T\mathbf{G}]\beta(\lambda, \theta) = \Phi^T\mathbf{y} + \theta\mathbf{G}^T\mathbf{c} \quad (3.2)$$

Note that $[\Phi^T \Phi + \lambda \mathbf{I}_m + \theta \mathbf{G}^T \mathbf{G}]$ is invertible for all λ and θ since Φ is $n \times m$ with rank m and $n > m$. Thus we can write the unique solution $\beta(\lambda, \theta)$ for 2.14 as:

$$\beta(\lambda, \theta) = [\Phi^T \Phi + \lambda \mathbf{I}_m + \theta \mathbf{G}^T \mathbf{G}]^{-1} \{\Phi^T \mathbf{y} + \theta \mathbf{G}^T \mathbf{c}\} \quad (3.3)$$

To improve computational stability, we can use the Sherman-Morrison-Woodbury formula to compute the matrix inverse in 3.3:

$$[\Phi^T \Phi + \lambda \mathbf{I}_m + \theta \mathbf{G}^T \mathbf{G}]^{-1} = \mathbf{R} - \mathbf{R} \mathbf{G}^T \left[\frac{1}{\theta} \mathbf{I}_q + \mathbf{G} \mathbf{R} \mathbf{G}^T \right]^{-1} \mathbf{G} \mathbf{R} \quad (3.4)$$

where

$$\mathbf{R} = (\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \quad (3.5)$$

Note also that:

$$\frac{\partial^2 f_{con}(\beta)}{\partial \beta \beta^T} = 2 [\Phi^T \Phi + \lambda \mathbf{I}_m + \theta \mathbf{G}^T \mathbf{G}] \quad (3.6)$$

Clearly $\frac{\partial^2 f_{con}(\beta)}{\partial \beta \beta^T}$ is positive definite and hence $\beta(\lambda, \theta)$ is the unique minimizer of 2.14.

4 The limit solution as $\theta \rightarrow \infty$

A simple way of controlling the "softness" of the constraint 2.13 is to vary the constraint controlling parameter θ in the objective function 2.14. In this section we study the properties of the soft constrained solution as $\theta \rightarrow \infty$.

Proposition 4.1. *Suppose $\beta(\lambda, \theta)$ is a solution to the soft constrained problem 2.14 then we have:*

1.

$$\lim_{\theta \rightarrow \infty} \mathbf{G} \beta(\lambda, \theta) = \mathbf{G} \lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = \mathbf{c} \quad (4.1)$$

2.

$$\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = (\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \{\Phi^T \mathbf{y} - \mathbf{G}^T \delta\} \quad (4.2)$$

where

$$\delta = [\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T]^{-1} \{\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \Phi^T \mathbf{y} - \mathbf{c}\} \quad (4.3)$$

Proof. Since $\beta(\lambda, \theta)$ is a solution to 2.14 it must satisfy the necessary conditions for optimality in 3.2. Pre-multiplying both sides of equation 3.2 by $\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1}$ we get:

$$[\mathbf{I}_q + \theta \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T] \mathbf{G}\beta(\lambda, \theta) = \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} [\Phi^T \mathbf{y} + \theta \mathbf{G}^T \mathbf{c}] \quad (4.4)$$

Since $(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1}$ is symmetric and positive definite, the $q \times q$ matrices $\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T$ and $[\mathbf{I}_q + \theta \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T]$ are both invertible. Therefore we get:

$$\mathbf{G}\beta(\lambda, \theta) = [\mathbf{I}_q + \theta \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T]^{-1} \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} [\Phi^T \mathbf{y} + \theta \mathbf{G}^T \mathbf{c}] \quad (4.5)$$

Re-arranging the right hand side we can write:

$$\mathbf{G}\beta(\lambda, \theta) = \left[\frac{\mathbf{I}_q}{\theta} + \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T \right]^{-1} \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \left[\frac{\Phi^T \mathbf{y}}{\theta} + \mathbf{G}^T \mathbf{c} \right] \quad (4.6)$$

Taking the limit as $\theta \rightarrow \infty$ on both sides and noting that $\frac{1}{\theta} \rightarrow 0$ as $\theta \rightarrow \infty$ we get:

$$\lim_{\theta \rightarrow \infty} \mathbf{G}\beta(\lambda, \theta) = \mathbf{G} \lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = [\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T]^{-1} \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} [\mathbf{G}^T \mathbf{c}] \quad (4.7)$$

Simplifying equation 4.7 we get:

$$\lim_{\theta \rightarrow \infty} \mathbf{G}\beta(\lambda, \theta) = \mathbf{G} \lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = \mathbf{c} \quad (4.8)$$

This proves part 1 of the proposition. Now re-arranging equation 3.2 we can write:

$$(\Phi^T \Phi + \lambda \mathbf{I}_m) \beta(\lambda, \theta) + \mathbf{G}^T \left[\underbrace{\theta}_{\text{tends to } \infty} \underbrace{\{\mathbf{G}\beta(\lambda, \theta) - \mathbf{c}\}}_{\text{tends to } 0} \right] = \Phi^T \mathbf{y} \quad (4.9)$$

In the above equation $\theta \rightarrow \infty$ and from 4.8 $\lim_{\theta \rightarrow \infty} \{\mathbf{G}\beta(\lambda, \theta) - \mathbf{c}\} \rightarrow \mathbf{0}$. Let

$$\lim_{\theta \rightarrow \infty} \theta \{\mathbf{G}\beta(\lambda, \theta) - \mathbf{c}\} = \delta \quad (4.10)$$

where δ is a $q \times 1$ vector. Taking $\lim_{\theta \rightarrow \infty}$ on both sides of 4.9 and using 4.10 we get:

$$(\Phi^T \Phi + \lambda \mathbf{I}_m) \lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) + \mathbf{G}^T \delta = \Phi^T \mathbf{y} \quad (4.11)$$

Pre-multiplying both sides of 4.11 by $\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1}$ we get:

$$\mathbf{G} \lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) + \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T \delta = \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \Phi^T \mathbf{y} \quad (4.12)$$

Using 4.8 and noting the non-singularity of $\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T$ we can solve for δ to get:

$$\delta = [\mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T]^{-1} \{ \mathbf{G}(\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \Phi^T \mathbf{y} - \mathbf{c} \} \quad (4.13)$$

With this value of δ we can use 4.11 to get:

$$\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = (\Phi^T \Phi + \lambda \mathbf{I}_m)^{-1} \{ \Phi^T \mathbf{y} - \mathbf{G}^T \delta \} \quad (4.14)$$

This proves part 2 of the proposition. \square \square

5 Equivalence to solution of KKT system

In this section, we will show that the solution to the hard constrained problem 2.15 as obtained by a direct solution of Karush-Kuhn-Tucker (KKT) optimality conditions coincides with the solution to the soft constrained problem as the constraint enforcing parameter $\theta \rightarrow \infty$.

Proposition 5.1. *Suppose $\beta(\lambda, \theta)$ is a solution to the soft constrained problem 2.14 and let $\beta^*(\lambda)$ be a solution to the hard constrained problem 2.15 then*

$$\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = \beta^*(\lambda) \quad (5.1)$$

Proof. First consider the hard constrained problem 2.15. The Lagrangian for this constrained problem is:

$$\mathcal{L}(\beta, \alpha) = f_{hcon}(\beta) - \alpha^T \{ \mathbf{G} \beta - \mathbf{c} \} \quad (5.2)$$

where α is a $q \times 1$ vector of Lagrange multipliers for the q constraints in 2.13. If $\beta^*(\lambda)$ is a solution to 2.15 and α^* the corresponding Lagrange multiplier then the KKT necessary conditions for optimality are:

$$\frac{\partial}{\partial \beta} \mathcal{L}(\beta, \alpha) |_{\beta^*(\lambda), \alpha^*} = \left[\frac{\partial}{\partial \beta} f_{hcon}(\beta) - \frac{\partial}{\partial \beta} \alpha^T \{ \mathbf{G} \beta - \mathbf{c} \} \right]_{\beta^*(\lambda), \alpha^*} = \mathbf{0} \quad (5.3)$$

$$\mathbf{G} \beta^*(\lambda) - \mathbf{c} = \mathbf{0} \quad (5.4)$$

Now from 2.15 we have:

$$\frac{\partial}{\partial \boldsymbol{\beta}} f_{hcon}(\boldsymbol{\beta})|_{\boldsymbol{\beta}^*(\lambda), \boldsymbol{\alpha}^*} = 2 (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m) \boldsymbol{\beta}^*(\lambda) - 2 \boldsymbol{\Phi}^T \mathbf{y} \quad (5.5)$$

Also

$$\frac{\partial}{\partial \boldsymbol{\beta}} \boldsymbol{\alpha}^T \{ \mathbf{G} \boldsymbol{\beta} - \mathbf{c} \} |_{\boldsymbol{\beta}^*(\lambda), \boldsymbol{\alpha}^*} = \mathbf{G}^T \boldsymbol{\alpha} |_{\boldsymbol{\beta}^*(\lambda), \boldsymbol{\alpha}^*} = \mathbf{G}^T \boldsymbol{\alpha}^* \quad (5.6)$$

Substituting 5.5 and 5.6 we can re-write the KKT optimality conditions as:

$$\begin{aligned} 2 (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m) \boldsymbol{\beta}^*(\lambda) - 2 \boldsymbol{\Phi}^T \mathbf{y} - \mathbf{G}^T \boldsymbol{\alpha}^* &= \mathbf{0} \\ \mathbf{G} \boldsymbol{\beta}^*(\lambda) - \mathbf{c} &= \mathbf{0} \end{aligned} \quad (5.7)$$

Pre-multiplying both sides of the first equation in 5.7 by $\frac{1}{2} \mathbf{G} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1}$ we get:

$$\mathbf{G} \boldsymbol{\beta}^*(\lambda) - \mathbf{G} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \boldsymbol{\Phi}^T \mathbf{y} - \mathbf{G} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T \left(\frac{\boldsymbol{\alpha}^*}{2} \right) = \mathbf{0} \quad (5.8)$$

Now $\mathbf{G} \boldsymbol{\beta}^*(\lambda) = \mathbf{c}$ from the second equation in 5.7. Noting the non-singularity of $(\mathbf{G} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T)$ we can solve for $(\frac{\boldsymbol{\alpha}^*}{2})$ to get:

$$\left(\frac{\boldsymbol{\alpha}^*}{2} \right) = [\mathbf{G} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \mathbf{G}^T]^{-1} \{ \mathbf{c} - \mathbf{G} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \boldsymbol{\Phi}^T \mathbf{y} \} \quad (5.9)$$

Comparing 5.9 and 4.3 from Proposition 4.1 it is clear that:

$$\left(\frac{\boldsymbol{\alpha}^*}{2} \right) = -\boldsymbol{\delta} \quad (5.10)$$

Solving for $\boldsymbol{\beta}^*(\lambda)$ from the first equation in 5.7 we get:

$$\boldsymbol{\beta}^*(\lambda) = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \left\{ \boldsymbol{\Phi}^T \mathbf{y} + \mathbf{G}^T \left(\frac{\boldsymbol{\alpha}^*}{2} \right) \right\} \quad (5.11)$$

From 5.10 and 5.11 we can write:

$$\boldsymbol{\beta}^*(\lambda) = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \{ \boldsymbol{\Phi}^T \mathbf{y} - \mathbf{G}^T \boldsymbol{\delta} \} \quad (5.12)$$

Comparing 5.12 and 4.14 we get:

$$\lim_{\theta \rightarrow \infty} \boldsymbol{\beta}(\lambda, \theta) = \boldsymbol{\beta}^*(\lambda) = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_m)^{-1} \{ \boldsymbol{\Phi}^T \mathbf{y} - \mathbf{G}^T \boldsymbol{\delta} \} \quad (5.13)$$

which proves the proposition. \square \square

6 Illustrative Experiment

To illustrate the results in Propositions 4.1 and 5.1 we performed a simple numerical experiment. We generated example data for \mathbf{y} and Φ as follows:

1. In MATLAB notation suppose $\mathbf{e} = \text{ones}(100, 1)$ and $\mathbf{u} = \text{linspace}(0, 1, 100)^T$. Suppose \mathbf{U} be a diagonal matrix with the elements of \mathbf{u} on the diagonal, in other words $U(i, i) = u_i$
2. Let $\Phi = [\mathbf{e}, \mathbf{u}, \mathbf{U}\mathbf{u}, \mathbf{U}^2\mathbf{u}, \mathbf{U}^3\mathbf{u}]$. Here \mathbf{u} is a linear basis function, $\mathbf{U}\mathbf{u}$ is a quadratic basis function, $\mathbf{U}^2\mathbf{u}$ is a cubic basis function and so on
3. We generated a random coefficient vector $\beta = \text{unifrnd}(-1, 1, 5, 1) \times 10$
4. Next we generate $\mathbf{y} = \Phi\beta$ and set $\lambda = 1$
5. We chose $\mathbf{G} = \text{unifrnd}(0, 1, 2, 5)$ and $\mathbf{c} = \text{unifrnd}(0, 1, 2, 1)$ to create a set of 2 linear constraints on β

Note that in this dummy example $n = 100$, $m = 5$ and $q = 2$. Next we chose a 170×1 vector of potential values for θ in 2.14 as:

$$\theta_{vec} = [\text{linspace}(0, 30, 100), \text{linspace}(30, 100, 50), \text{linspace}(100, 1000, 20)]^T$$

For each i from 1 to 170 we set $\theta = \theta_{vec}(i)$ and solve the soft constrained problem 2.14 using 3.3 to get $\beta(\lambda, \theta)$. For the same θ we also solve the hard constrained problem 2.15 using 5.13 and 4.3 to get $\beta^*(\lambda)$. By Proposition 5.1 we have $\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = \beta^*(\lambda)$. To verify this numerically we simply plot the components of $\beta(\lambda, \theta)$ on the x -axis and the corresponding components of $\beta^*(\lambda)$ on the y -axis. As $\beta(\lambda, \theta)$ gets closer and closer to $\beta^*(\lambda)$ with increasing θ , the plotted points should get closer to the $y = x$ line and eventually lie on the $y = x$ line. Figure 1 shows a plot of $\beta(\lambda, \theta)$ versus $\beta^*(\lambda)$ for 4 values $\theta = 0, 4, 20, 1000$. It can be seen that as θ increases the components of $\beta(\lambda, \theta)$ converge to the components of $\beta^*(\lambda)$.

7 Discussion

In this work we considered least squares problems with smoothness regularization and enforcement of constraints in the soft form by setting up the optimization problem 2.14. First, we solve 2.14 to get a solution $\beta(\lambda, \theta)$

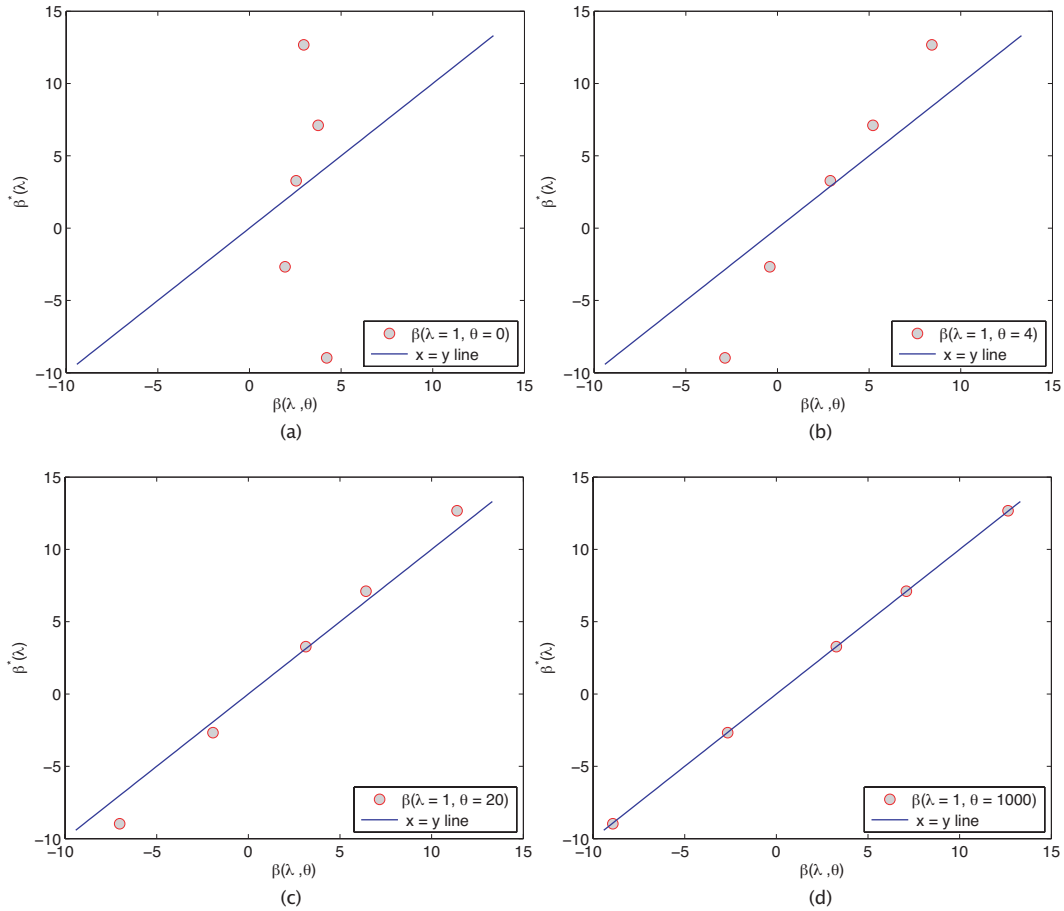


Figure 1: Figure shows the behavior of regularized and soft constrained solutions $\beta(\lambda, \theta)$ for various values of θ for a fixed $\lambda = 1$. The x -axis shows components of $\beta(\lambda, \theta)$. The corresponding components of the hard constrained solution $\beta^*(\lambda)$ are shown on the y -axis. It can be seen that as $\theta \rightarrow \infty$ the agreement between $\beta(\lambda, \theta)$ and $\beta^*(\lambda)$ increases.

that is parameterized by λ and θ . Next, we study the behavior of this solution as $\theta \rightarrow \infty$ in Proposition 4.1 and show in Proposition 5.1 that in fact $\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta) = \beta^*(\lambda)$ where $\beta^*(\lambda)$ is the solution to the hard constrained optimization problem 2.15. We also illustrate this convergence phenomenon via a simple illustrative experiment in Section 6.

Enforcement of constraints in a soft manner as in 2.14 is similar to the quadratic penalty method in constrained optimization [2]. Results related to limiting behavior such as 5.1 have also been proven under more general conditions in [2] given the assumption of convexity of the intermediate quadratic penalty problems. It is worth noting that the limiting process $\lim_{\theta \rightarrow \infty} \beta(\lambda, \theta)$ could be difficult to carry out numerically because of the ill conditioning inherent in the problem as $\theta \rightarrow \infty$ [2]. However, we have circumvented this difficulty in the present case by using the Sherman-Morrison-Woodbury analytical matrix inversion in 3.4.

8 Conclusion

Solution to problems such as 2.14 could find many practical applications since they not only impose smoothness via the parameter λ on the inferred function but also enable the modeler to impose soft constraints via the parameter θ . By varying λ and θ , a family of solutions with various properties is obtained. It is re-assuring to know that increasing θ will eventually lead to the fully constrained solution corresponding to the problem 2.15.

References

- [1] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, 2001.
- [2] J. Nocedal and S.J. Wright. *Numerical Optimization*. New York:Springer, 1999.