

51 REGULARIZATION

We discussed the least-squares problem in the last chapter, which uses a collection of data points $\{x(n), y_n\}$ to determine an optimal parameter w^* by minimizing an empirical quadratic risk of the form:

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ P(w) \triangleq \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \quad (51.1a)$$

where each y_n is M -dimensional and each $x(n)$ is a scalar. The solution is determined by solving the normal equations:

$$H^\top H w^* = H^\top d, \quad (\text{normal equations}) \quad (51.1b)$$

where the quantities $d \in \mathbb{R}^{N \times 1}$ and $H \in \mathbb{R}^{N \times M}$ collect the data:

$$H \triangleq \begin{bmatrix} y_0^\top \\ y_1^\top \\ y_2^\top \\ \vdots \\ y_{N-1}^\top \end{bmatrix}, \quad d \triangleq \begin{bmatrix} x(0) \\ x(1) \\ x(2) \\ \vdots \\ x(N-1) \end{bmatrix} \quad (51.1c)$$

The normal equations (51.1b) may have a unique solution or infinitely many solutions. They may also be ill-conditioned meaning that slight perturbations to the data $\{d, H\}$ can lead to large changes in the solution w^* ; this usually occurs when the matrix H is ill-conditioned. In this chapter, we will use the least-squares formulation as a guiding example to illustrate three types of challenges that arise in data-driven learning methods pertaining to **(a)** non-uniqueness of solutions, **(b)** ill-conditioning, and **(c)** the undesirable possibility of over-fitting. We will then explain that regularization is a useful tool to alleviate these challenges. We will also explain how regularization enables the designer to promote preference for certain solutions such as favoring solutions with small norms or sparse structure. We will motivate the main ideas by using the least-squares formulation due to its mathematical tractability. Subsequently, we will extend the discussion more general empirical risks, other than the least-squares case, which will arise in later chapters when we deal with logistic regression, support vector machines, kernel machines, neural networks, and other learning methods.

51.1 THREE CHALLENGES

In learning problems, we make a distinction between *training* data and *test* data. The data $\{x(n), y_n\}$ used to solve the least-squares problem (51.1a) are referred to as *training data*. Once a solution w^* is determined, the value of the risk function at the solution is called the *training error*:

$$\text{training error} \triangleq \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w^*)^2 = \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - \hat{x}(n))^2 \quad (51.2)$$

where $\hat{x}(n) = y_n^\top w^*$ denotes the prediction for $x(n)$. In this way, the training error is measuring how well the least-squares solution performs on the training data. In general, the training error will be small because the solution w^* is purposefully determined to minimize it.

In most learning applications, however, the main purpose for learning w^* is to employ it to perform prediction on future data that *were not* part of the training phase. For this reason, it is customary to assess performance on a separate collection of T test data points denoted by $\{x(t), y_t\}$, and which are assumed to arise from the same underlying distribution $f_{\mathbf{x}, \mathbf{y}}(x, y)$ as the training data. The corresponding testing error is defined by

$$\text{testing error} \triangleq \frac{1}{T} \sum_{t=0}^{T-1} (x(t) - y_t^\top w^*)^2 = \frac{1}{T} \sum_{t=0}^{T-1} (x(t) - \hat{x}(t))^2 \quad (51.3)$$

where $\hat{x}(t) = y_t^\top w^*$ denotes the prediction for $x(t)$. In general, the testing error will be larger than the training error but we desire the gap between them to be small. Learning algorithms that lead to small error gaps are said to *generalize* well, namely, they are able to extend their good performance on training data to the test data as well. We will discuss generalization and training and testing errors in greater detail in future chapters, especially in the context of classification problems. Here, we are using the least-squares problem to motivate the concepts.

Difficulties

We already know that the normal equations (51.1b) are consistent, meaning that a solution w^* always exists. The solution is either unique when H has full column rank, in which case it is given by

$$w^* = (H^\top H)^{-1} H^\top d, \quad (H \text{ has full column rank}) \quad (51.4)$$

or there are infinitely many solutions differing by vectors in $\mathcal{N}(H)$. Some challenges arise in both scenarios, which lead to complications when solving inference problems:

(a) (Non-uniqueness). When infinitely many solutions exist, the training error

will not change regardless of which solution we pick. This is because all valid solutions w^* differ by vectors in the null space of H and, therefore, if w_1^* and w_2^* are two valid solutions then

$$w_2^* = w_1^* + p, \quad \text{for some } p \in \mathcal{N}(H) \quad (51.5)$$

In this case, the predictions $\hat{x}(n)$ for the training signals will remain unchanged under w_1^* or w_2^* since $Hp = 0$ and, hence, $y_n^\top p = 0$ for any of the observation vectors in the training set so that

$$\hat{x}(n) = y_n^\top w_2^* = y_n^\top w_1^* \quad (51.6)$$

It follows that the training error remains invariant. However, the testing error will be sensitive to which solution we select because the *test* observations $\{y_t\}$ need not be orthogonal anymore to the nullspace of H . We explain in the sequel that ℓ_2 -regularization forces a unique solution w^* and removes this ambiguity.

- (b) (**Overfitting**) Infinitely many solutions w^* can exist even when $N \geq M$, i.e., even when we have more observations than unknown entries. This occurs when the columns of H are linearly dependent and gives rise to a second challenge. Recall that the least-squares problem is approximating d by $\hat{d} = Hw^*$. When the columns of H are linearly dependent, some of its columns can be removed to obtain a full-rank lower-dimensional matrix, $H' \in \mathbb{R}^{N \times M'}$ with $M' < M$. This new matrix spans the same column space as H :

$$\mathcal{R}(H') = \mathcal{R}(H) \quad (51.7)$$

We can then solve an equivalent least-squares problem involving $\{d, H'\}$ instead of $\{d, H\}$ to obtain the same projection \hat{d} by using a smaller-size solution $(w')^*$ of dimension M' . We thus see that the rank-deficiency of H amounts to using a more complex model w (i.e., of higher dimensions) than is necessary to approximate d . This issue is a manifestation of the problem of *overfitting*, which we will discuss in greater detail in later chapters. Overfitting amounts to using more complex models than necessary and it also degrades performance on test data.

Rank-deficiency of H also arises when $N < M$ (i.e., when H has more columns than rows). One way to deal with this problem is to collect more data (i.e., to use a larger N). A second way is to perform dimensionality reduction and reduce the size of the observation vectors. We will discuss techniques for dimensionality reduction in later chapters, including the principal component analysis (PCA) method and the Fisher discriminant analysis (FDA) method. A third way is to employ regularization. For example, we will explain further ahead that ℓ_1 -regularization automatically selects a subset of the columns of H to compute w^* .

- (c) (**Ill-conditioning**) Difficulties can arise even when the normal equations have a unique solution w^* but the data matrix H is ill-conditioned (i.e., has a large

condition number). In this case, small changes in the data $\{d, H\}$ can lead to large changes in the solution w^* and affect the inference conclusion and testing error — see Prob. 51.2 for a numerical example. One leading cause for ill-conditioning is when the entries within the observation vectors are not normalized properly so that some entries are disproportionately larger by some orders of magnitude than other entries. Such large discrepancies can distort the operation of a learning algorithm, including the least-squares solution, by giving more relevance or attention to larger entries in the observation vector over other entries. One way to deal with ill-conditioning is therefore to *scale* the observation vectors so that their entries assume values within some uniform range. The next example explains how scaling can be performed. A second way is to employ regularization. In particular, we will see that ℓ_2 -regularization reduces the effect of ill-conditioning.

Example 51.1 (Normalization of observation vectors) It is common practice to center the training data around their sample means, as was already suggested by the discussion in Sec. 29.2. We can take this step further and normalize the entries of the observation vectors to have unit-variance as well. Specifically, the first step is to compute the sample mean vector:

$$\bar{y} \triangleq \frac{1}{N} \sum_{n=0}^{N-1} y_n \quad (51.8a)$$

and to use it to center all observation vectors by replacing them by

$$y_{n,c} \triangleq y_n - \bar{y} \quad (51.8b)$$

where, for clarity, we are adding the subscript “*c*” to refer to centered variables. If we denote the individual entries of $\{\bar{y}, y_n\}$ by $\{\bar{y}(m), y(m), m = 1, 2, \dots, M\}$, then centering amounts to replacing the individual entries by

$$y_{n,c}(m) \triangleq y_n(m) - \bar{y}(m) \quad (51.8c)$$

The second step in the normalization process is to evaluate the (unbiased) sample variance for each of these centered entries, namely,

$$\hat{\sigma}_m^2 \triangleq \frac{1}{N-1} \sum_{n=0}^{N-1} y_{n,c}^2(m), \quad m = 1, 2, \dots, M \quad (51.9a)$$

and to scale $y_{n,c}(m)$ by the corresponding standard deviation to get

$$y_{n,p}(m) \triangleq y_{n,c}(m) / \hat{\sigma}_m, \quad m = 1, 2, \dots, M \quad (51.9b)$$

where we are now using the subscript “*p*.” In this way, we start from an observation vector y_n and replace it by the normalized vector $y_{n,p}$, where all entries of $y_{n,p}$ are centered with zero mean and unit variance:

$\{y_n\}$	$\xrightarrow{\text{remove sample mean}}$	$\{y_{n,c}\}$	$\xrightarrow{\text{normalize variance}}$	$\{y_{n,p}\}$
-----------	---	---------------	---	---------------

(51.10)

A second method to normalize the observation vectors $\{y_n\}$ is as follows. We first identify the smallest and largest entry values within the given dataset:

$$y_{\min} \triangleq \min_{n,m} y_n(m) \quad (51.11a)$$

$$y_{\max} \triangleq \max_{n,m} y_n(m) \quad (51.11b)$$

$$\Delta = y_{\max} - y_{\min} \quad (51.11c)$$

and then scale all entries in the following manner, for each n and m :

$$y_{n,s}(m) \triangleq \frac{y_n(m) - y_{\min}}{\Delta} \quad (51.12)$$

In this way, each scaled entry $y_{n,s}(m)$ will assume values within the range $[0, 1]$. We can subsequently center the means of these entries at zero by computing

$$y_{n,p} \triangleq y_{n,s} - \bar{y}_{n,s}, \quad \text{where} \quad \bar{y}_{n,s} = \frac{1}{N} \sum_{n=0}^{N-1} y_{n,s} \quad (51.13)$$

Here again, we start from a given observation vector y_n and replace it by $y_{n,p}$, where all entries lie within the range $[-1, 1]$:

$$\boxed{\begin{array}{ccccc} & \text{normalize} & & \text{remove} & \\ & \text{range} & & \text{sample} & \\ \{y_n\} & \xrightarrow{\quad} & \{y_{n,s}\} & \xrightarrow{\quad} & \{y_{n,p}\} \end{array}} \quad (51.14)$$

Regardless of which normalization procedure is used, we will assume that the given observation vectors $\{y_n\}$ have already gone through this process and will continue to use the notation y_n rather than switch to $y_{n,p}$ for simplicity.

51.2 ℓ_2 -REGULARIZATION

One useful technique to avoid the challenges of non-uniqueness of solutions, overfitting, and ill-conditioning is to employ *regularization* (also called *shrinkage* in the statistics literature). The technique penalizes some norm of the parameter w in order to favor solutions with desirable properties based on some prior knowledge (such as sparse solutions or solutions with small Euclidean norm). We say that regularization incorporates a form of *inductive bias* in that it biases the solution away from the unregularized case by incorporating some prior information. This is attained by adding an explicit convex penalty term to the original risk function such as

$$q(w) = \begin{cases} \rho \|w\|^2, & (\ell_2\text{-regularization}) \\ \alpha \|w\|_1, & (\ell_1\text{-regularization}) \\ \alpha \|w\|_1 + \rho \|w\|^2, & (\text{elastic-net regularization}) \\ \beta \|w\|_0, & (\ell_0\text{-regularization}) \end{cases} \quad (51.15)$$

where (α, β, ρ) are nonnegative parameters, and where $\|w\|_0$ is a pseudo-norm that counts the number of nonzero elements in w . We will focus on the first

three choices due to their mathematical tractability. One can also consider other vector norms, such as the p -th norm, $\|w\|_p$ for $p < 1$ or $p = \infty$. Regularization will generally have a limited effect on the *training error* of an algorithm, but will improve the *generalization* ability of the algorithm by improving its performance on test data for the reasons explained in the sequel. We consider first the case of ℓ_2 -regularization, also called *ridge regression*, where the penalty term is quadratic in w .

51.2.1 Ridge Regression

In ridge regression, we replace the empirical risk (51.1a) by the regularized version:

$$w_{\text{reg}}^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ P_{\text{reg}}(w) \triangleq \rho \|w\|^2 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \quad (51.16)$$

where $\rho > 0$ is the regularization factor; its value may or may not depend on N . In general, the value of ρ is independent of N .

Observe that, for the purposes of this chapter, we are adding a subscript “reg” to $(w_{\text{reg}}^*, P_{\text{reg}}(w))$ to distinguish them from the unregularized versions $(w^*, P(w))$. This is because we will be comparing both risks and their minimizers throughout this chapter. In future chapters, however, where we will be working almost exclusively with regularized risks, we will revert to the traditional notation $(w^*, P(w))$ without the “reg” subscript for simplicity. Before explaining how ridge regression addresses the aforementioned challenges, we revisit Example 50.1 and show how the regularized empirical risk (51.16) can be motivated as the solution to a maximum a-posteriori (MAP) inference problem.

Example 51.2 (Interpretation in terms of a Gaussian prior on the model) Assume we collect N independent and identically-distributed observations $\{\mathbf{x}(n), y_n\}$, for $0 \leq n \leq N-1$. Assume also that these observations satisfy the same linear model (50.20), namely,

$$\mathbf{x}(n) = y_n^\top w + \mathbf{v}(n) \quad (51.17)$$

for some unknown $w \in \mathbb{R}^M$, and where $\mathbf{v}(n)$ is a white Gaussian noise process with zero mean and variance σ_v^2 . In the earlier Example 50.1, the model w was treated as an *unknown constant* and a maximum-likelihood formulation was used to estimate it; thus leading to the standard least-squares problem. Here, we will instead model w as a realization for some random variable \mathbf{w} that is Gaussian-distributed with zero mean and covariance matrix $R_w = \sigma_w^2 I_M$, i.e.,

$$f_{\mathbf{w}}(w) = \frac{1}{\sqrt{(2\pi\sigma_w^2)^M}} \exp\left\{-\frac{1}{2\sigma_w^2} \|w\|^2\right\} \quad (51.18)$$

Once w is selected from this distribution, then all observations $\{\mathbf{x}(n)\}$ are generated by this *same* w from knowledge of $\{x(n), y_n\}$. We are again interested in estimating w . Using Bayes rule (3.39), we assess the conditional probability distribution of the model

given the observations as follows:

$$\begin{aligned}
& f_{\mathbf{w}|\mathbf{x},\mathbf{y}}(w|\{x(n), y_n\}) \\
& \propto f_{\mathbf{x},\mathbf{y}|\mathbf{w}}(\{x(n), y_n\} | w) f_{\mathbf{w}}(w) \\
& = \left\{ \prod_{n=0}^{N-1} f_v(x(n) - y_n^\top w) \right\} f_{\mathbf{w}}(w) \\
& \propto \left\{ \prod_{n=0}^{N-1} \exp\left\{-\frac{1}{2\sigma_v^2}(x(n) - y_n^\top w)^2\right\} \right\} \times \exp\left\{-\frac{1}{2\sigma_w^2}\|w\|^2\right\} \\
& = \exp\left\{-\frac{1}{2\sigma_w^2}\|w\|^2 - \frac{1}{2\sigma_v^2} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2\right\} \tag{51.19}
\end{aligned}$$

where the first and third lines replace the equality sign by proportionality constants. Consequently, we can now formulate a *maximum a-posteriori* (MAP) estimation problem to recover w , which amounts to seeking the value of w that maximizes the above conditional density function:

$$\begin{aligned}
w_{\text{reg}}^* & \triangleq \operatorname{argmax}_{w \in \mathbb{R}^M} f_{\mathbf{w}|\mathbf{x},\mathbf{y}}(w|\{x(n), y_n\}) \\
& = \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \frac{1}{2\sigma_w^2}\|w\|^2 + \frac{1}{2\sigma_v^2} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \\
& = \operatorname{argmin}_{w \in \mathbb{R}^M} \frac{N}{2\sigma_v^2} \left\{ \frac{2\sigma_v^2}{2N\sigma_w^2}\|w\|^2 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \\
& = \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \rho\|w\|^2 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \tag{51.20}
\end{aligned}$$

where we introduced $\rho = \sigma_v^2/N\sigma_w^2$. We therefore recover the regularized empirical risk (51.16). This argument shows that ℓ_2 -regularization helps ensure that the solution w_{reg}^* is consistent with a prior Gaussian model on the distribution of \mathbf{w} .

We now explain how ridge regression promotes solutions with smaller Euclidean norm and alleviates the challenges of ill-conditioning, over-fitting, and non-uniqueness of solutions.

Resolving non-uniqueness

Differentiating $P_{\text{reg}}(w)$ in (51.16) with respect to w , we find that the solution is unique and given by

$$w_{\text{reg}}^* = (\rho N I_M + H^\top H)^{-1} H^\top d \tag{51.21}$$

where the matrix $\rho N I_M + H^\top H$ is always invertible due to the positive term, $\rho N I_M > 0$ and independent of whether H is rank-deficient or not.

Promoting smaller solutions

It is seen from the regularized risk (51.16) that larger values for ρ favor solutions w_{reg}^* with smaller Euclidean norm than would result when $\rho = 0$. This is because

the objective is to minimize the aggregate risk, and the first term is influenced by $\rho\|w\|^2$. This property can be established more formally as follows (see Prob. 51.3 for an alternative argument). Using the unregularized risk $P(w)$, and since w_{reg}^* minimizes the regularized risk, we have

$$\begin{aligned} \rho\|w_{\text{reg}}^*\|^2 + P(w_{\text{reg}}^*) &\leq \rho\|w^*\|^2 + P(w^*) \\ \implies \rho\|w_{\text{reg}}^*\|^2 - \rho\|w^*\|^2 &\leq P(w^*) - P(w_{\text{reg}}^*) \\ \stackrel{(a)}{\implies} \rho\|w_{\text{reg}}^*\|^2 - \rho\|w^*\|^2 &\leq 0 \end{aligned} \quad (51.22)$$

where step (a) is because w^* minimizes the unregularized risk, $P(w)$. It follows that $\|w_{\text{reg}}^*\|^2 \leq \|w^*\|^2$. Actually, strict inequality holds because $P(w^*)$ is strictly smaller than $P(w_{\text{reg}}^*)$. Since, otherwise, for $P(w)$ to assume the same value at both (w^*, w_{reg}^*) , it would mean that w_{reg}^* must be a minimizer for $P(w)$ as well. In that case, both (w^*, w_{reg}^*) must satisfy the *same* normal equations, namely,

$$H^T H w^* = H^T d, \quad H^T H w_{\text{reg}}^* = H^T d \quad (51.23)$$

But since w_{reg}^* satisfies (51.21), i.e., $(\rho N I_M + H^T H) w_{\text{reg}}^* = H^T d$, we conclude that $w_{\text{reg}}^* = 0$. But this is not possible unless $H^T d = 0$. Absent this condition, we conclude that

$$\|w_{\text{reg}}^*\|^2 < \|w^*\|^2 \quad (51.24)$$

This proves that the norm of the regularized solution, w_{reg}^* , shrinks in comparison to the norm of the original solution, w^* . This property is referred to as *shrinkage*. We will encounter it in other regularization formulations as well.

Countering ill-conditioning

Regularization also counters the effect of ill-conditioning, i.e., the sensitivity of the solution w^* to small variations in the data $\{x(n), y_n\}$. Note that the condition number of the new coefficient matrix is given by

$$\kappa(\rho N I_M + H^T H) \triangleq \frac{\rho N + \lambda_{\max}(H^T H)}{\rho N + \lambda_{\min}(H^T H)} = \frac{\rho N + \sigma_{\max}^2(H)}{\rho N + \sigma_{\min}^2(H)} \quad (51.25)$$

in terms of the largest and smallest singular values of H . If the value of ρN is large enough in comparison to the singular-value spread of H , then the ratio on the right-hand side approaches one and the matrix $\rho N I_M + H^T H$ becomes (very) well conditioned.

Countering overfitting

By promoting solutions w_{reg}^* with smaller Euclidean norm, regularization helps alleviate the danger of overfitting because it searches for the solution over a reduced region in space. This can be shown more formally by verifying that minimizing a regularized least-squares problem of the form (51.16) is equivalent to solving a *constrained* optimization problem of the following form:

$$w_{\text{reg}}^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^T w)^2 \right\}, \quad \text{subject to } \rho\|w\|^2 \leq \tau \quad (51.26)$$

for some $\tau > 0$. The equivalence between problems (51.16) and (51.26) is established algebraically in Appendix 51.A by using the Lagrange and KKT multiplier arguments from Sec. 9.1. This equivalent characterization shows that regularization reduces the search space for w to the spherical region $\|w\|^2 \leq \tau/\rho$ instead of searching over the entire space $w \in \mathbb{R}^M$. Some care is needed in selecting ρ (or τ): large values for ρ (or small τ) can have the opposite effect and constrain the search region excessively, thus leading to the possibility of *underfitting* (i.e., to the use of simpler models than is actually necessary to fit the data well). These remarks show that there is a compromise in setting the value of ρ : small ρ does not perform effective regularization and large ρ can cause underfitting.

Biased risk values

Although regularization is effective in countering ill-conditioning and overfitting, there is a price to pay. This is because regularization biases the least attainable risk (i.e., the training error), which becomes larger than in the unregularized case. To see this, consider again the solutions w^* and w_{reg}^* to the standard and regularized least-squares problems. Evaluating the risk functions at the respective minimizers and subtracting them we get, after some algebra — see Prob. 51.4:

$$P_{\text{reg}}(w_{\text{reg}}^*) - P(w^*) = \rho (w^*)^T w_{\text{reg}}^* > 0 \quad (51.27)$$

from which we conclude that $P_{\text{reg}}(w_{\text{reg}}^*) > P(w^*)$, and that the bias increases with ρ .

Example 51.3 (QR solution method) Determination of the ℓ_2 -regularized solution (51.21) requires that we compute the matrix product $H^T H$ and invert the matrix $\rho N I_M + H^T H$. We explained earlier in Prob. 50.5 that squaring matrix entries through the product $H^T H$ can lead to a loss in numerical precision for small entries; it can also lead to overflow for large entries. A more stable numerical procedure for determining w_{reg}^* can be motivated by using the QR decomposition. We construct the extended quantities:

$$H^e \triangleq \begin{bmatrix} H \\ \sqrt{\rho N} I_M \end{bmatrix}, \quad \text{of size } (N+M) \times M \quad (51.28)$$

$$d^e \triangleq \begin{bmatrix} d \\ 0_{M \times 1} \end{bmatrix}, \quad \text{of size } (N+M) \times 1 \quad (51.29)$$

and introduce the QR decomposition:

$$H^e = Q \begin{bmatrix} R \\ 0 \end{bmatrix} \quad (51.30)$$

where

$$R: (M \times M), \quad Q: (N+M) \times (N+M), \quad Q Q^T = Q^T Q = I \quad (51.31)$$

and R is upper triangular. We apply the orthogonal transformation Q^T to d^e and denote the resulting entries by

$$Q^T d^e = \begin{bmatrix} \bar{d} \\ \times \end{bmatrix}, \quad \bar{d}: (N \times 1) \quad (51.32)$$

where \times refers to irrelevant entries. Then, note from (51.28) that

$$(H^e)^\top H^e = \rho N I_M + H^\top H \quad (51.33)$$

while from (51.30)

$$(H^e)^\top H^e = \begin{bmatrix} R^\top & 0 \end{bmatrix} Q^\top Q \begin{bmatrix} R \\ 0 \end{bmatrix} = R^\top R \quad (51.34)$$

It then follows that

$$\begin{aligned} w_{\text{reg}}^* &= (\rho N I_M + H^\top H)^{-1} H^\top d \\ &= \left((H^e)^\top H^e \right)^{-1} (H^e)^\top \begin{bmatrix} d \\ 0 \end{bmatrix} \end{aligned} \quad (51.35)$$

$$\begin{aligned} &= R^{-1} (R^\top)^{-1} (H^e)^\top d^e \\ &\stackrel{(51.30)}{=} R^{-1} (R^\top)^{-1} \begin{bmatrix} R^\top & 0 \end{bmatrix} Q^\top d^e \\ &\stackrel{(51.32)}{=} R^{-1} \begin{bmatrix} I_M & 0 \end{bmatrix} \begin{bmatrix} \bar{d} \\ \times \end{bmatrix} \\ &= R^{-1} \bar{d} \end{aligned} \quad (51.36)$$

We therefore arrive at the QR procedure listed in (51.37) for determining the ℓ_2 -regularized solution, which involves solving a triangular system of equations.

QR method for minimizing ℓ_2 -regularized least-squares risk (51.16).

$$\begin{aligned} &\text{given } \rho > 0 \text{ and data } d = \text{col}\{x(n)\}, H = \text{blkrow}\{y_n^\top\}; \\ &\text{construct } H^e = \begin{bmatrix} H \\ \sqrt{\rho N} I_M \end{bmatrix} \text{ and } d^e = \begin{bmatrix} d \\ 0_M \end{bmatrix}; \\ &\text{perform the QR decomposition } H^e = Q \begin{bmatrix} R \\ 0 \end{bmatrix}; \\ &\text{apply } Q^\top \text{ to } d^e \text{ and find } Q^\top d = \begin{bmatrix} \bar{d} \\ \times \end{bmatrix}; \\ &\text{solve the triangular system of equations } R w_{\text{reg}}^* = \bar{d}. \end{aligned} \quad (51.37)$$

51.3 ℓ_1 -REGULARIZATION

In ℓ_1 -regularization, we replace the empirical risk (51.1a) by the regularized version:

$$w_{\text{reg}}^* \triangleq \underset{w \in \mathbb{R}^M}{\text{argmin}} \left\{ P_{\text{reg}}(w) \triangleq \alpha \|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \quad (51.38)$$

in terms of the ℓ_1 -norm of w (i.e., the sum of its absolute entries), and where $\alpha > 0$ is the regularization factor; its value may or may not depend on N . In general, the value of α is independent of N . The variant with elastic-net

regularization solves instead

$$w_{\text{reg}}^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ P_{\text{reg}}(w) = \alpha \|w\|_1 + \rho \|w\|^2 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \quad (51.39)$$

with both $\alpha > 0$ and $\rho > 0$. We will discover in this section that ℓ_1 -regularization leads to a sparse solution w_{reg}^* , i.e., to a solution with a few nonzero entries. In this way, for any observation vector y , the inner product calculation $\hat{x} = y^\top w_{\text{reg}}^*$ ends up using only a few select entries from y due to the sparsity of w_{reg}^* . This means that ℓ_1 -regularization performs a form of “dimensionality reduction.” In particular, when some entries in y are correlated or redundant, the ℓ_1 -solution will rely on one of them and ignore the others. Elastic-net regularization, on the other hand, inherits useful features from both ℓ_2 and ℓ_1 -regularization. For example, it can handle situations involving more unknowns than measurements ($M > N$), and it also performs entry selection albeit in a less dramatic fashion than ℓ_1 -regularization.

The following derivation extends Example 51.2 and provides a similar MAP interpretation for the ℓ_1 -regularized empirical risk function (51.38).

Example 51.4 (Interpretation in terms of a Laplacian prior on the model) We collect N independent and identically-distributed observations $\{\mathbf{x}(n), y_n\}$, for $0 \leq n \leq N-1$, and assume that they satisfy the same linear model (51.17). The main difference is that we now assume that w is a realization of a random vector \mathbf{w} whose entries $\{w_m\}$ are independent of each other and arise from a Laplace distribution with zero mean and variance σ_w^2 :

$$f_{\mathbf{w}_m}(w_m) = \frac{1}{\sqrt{2}\sigma_w} \exp\{-\sqrt{2}|w_m|/\sigma_w\} \quad (51.40)$$

We also assume that all observations $\{\mathbf{x}(n)\}$ are generated by the *same* realization w . We are again interested in estimating w . Using Bayes rule (3.39), we assess the conditional probability distribution of the model given the observations as follows:

$$\begin{aligned} & f_{\mathbf{w}|\mathbf{x},\mathbf{y}}(w|\{\mathbf{x}(n), y_n\}) \\ & \propto f_{\mathbf{x},\mathbf{y}|\mathbf{w}}(\{\mathbf{x}(n), y_n\} | w) f_{\mathbf{w}}(w) \\ & = \left\{ \prod_{n=0}^{N-1} f_v(x(n) - y_n^\top w) \right\} f_{\mathbf{w}}(w) \\ & \propto \prod_{n=0}^{N-1} \exp\left\{-\frac{1}{2\sigma_v^2} (x(n) - y_n^\top w)^2\right\} \times \prod_{m=1}^M \exp\{-\sqrt{2}|w_m|/\sigma_w\} \\ & = \exp\left\{-\frac{\sqrt{2}}{\sigma_w} \|w\|_1 - \frac{1}{2\sigma_v^2} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2\right\} \end{aligned} \quad (51.41)$$

where the first and third lines replace the equality sign by proportionality constants. Consequently, we can now formulate a *maximum a-posteriori* (MAP) estimation prob-

lem to recover w by maximizing the above conditional density function as follows:

$$\begin{aligned}
 w_{\text{reg}}^* &\triangleq \operatorname{argmax}_{w \in \mathbb{R}^M} f_{w|x,h}(w|\{x(n), y_n\}) \\
 &= \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \frac{\sqrt{2}}{\sigma_w} \|w\|_1 + \frac{1}{2\sigma_v^2} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \\
 &= \operatorname{argmin}_{w \in \mathbb{R}^M} \frac{N}{2\sigma_v^2} \left\{ \frac{2\sqrt{2}\sigma_v^2}{N\sigma_w} \|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \\
 &= \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \alpha \|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \tag{51.42}
 \end{aligned}$$

where we introduced $\alpha = 2\sqrt{2}\sigma_v^2/N\sigma_w$. We therefore recover the regularized empirical risk (51.38) with $q(w) = \alpha\|w\|_1$. This argument shows that ℓ_1 -regularization helps ensure that the solution w_{reg}^* is consistent with a prior Laplacian model on the distribution of w .

We now explain how ℓ_1 -regularization (or its extension in terms of elastic-net regularization) promotes solutions with smaller norm and alleviates the challenges of ill-conditioning, overfitting, and non-uniqueness of solutions.

Resolving non-uniqueness

The penalty term $\alpha\|w\|_1$ is only convex. The regularized risk function will have a unique minimizer w_{reg}^* if the unregularized risk $P(w)$ happens to be strictly or strongly convex. For the least-squares case, the unregularized risk is strongly convex when $H^\top H > 0$. More generally, if this condition does not hold, then elastic-net regularization can be used and it will ensure a unique minimizer w_{reg}^* because the resulting regularized risk in that case will become strongly-convex regardless of whether H is rank-deficient or not.

Promoting smaller solutions

It is seen from the regularized risk in (51.38)–(51.39) that larger values for α or ρ favor solutions w_{reg}^* with smaller norms than would result when $\alpha = \rho = 0$. This is because the objective is to minimize the aggregate risk, and the regularization factors are influenced by $\alpha\|w\|_1$ and $\rho\|w\|^2$. This conclusion can be established more formally. If we set $q(w) = \alpha\|w\|_1$ for ℓ_1 -regularization or $q(w) = \alpha\|w\|_1 + \rho\|w\|^2$ for elastic-net regularization, then it follows from the general result in Appendix 51.A that the following shrinkage property holds:

$$q(w_{\text{reg}}^*) \leq q(w^*) \tag{51.43}$$

The result in the appendix holds for more general convex risks, $P(w)$, and is not limited to least-squares risks. It also holds for general convex regularization factors, $q(w)$, than ℓ_1 or elastic-net regularization. In other words, result (51.43) extends (51.24) to general convex risks and penalty terms.

Countering overfitting

Both ℓ_1 and elastic-net regularization help alleviate the danger of overfitting because they can also be shown to search for their solutions over reduced regions in space. This can be established more formally by verifying that minimizing a regularized least-squares problem of either forms (51.38)–(51.39) is equivalent to solving a *constrained* optimization problem of the following form:

$$w_{\text{reg}}^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^T w)^2 \right\}, \quad \text{subject to } q(w) \leq \tau \quad (51.44)$$

for some $\tau > 0$ and using the appropriate regularization factor: $q(w) = \alpha \|w\|_1$ for ℓ_1 -regularization and $q(w) = \alpha \|w\|_1 + \rho \|w\|^2$ for elastic-net regularization. The equivalence between problems (51.38)–(51.39) and (51.44) is again established algebraically in Appendix 51.A by using the KKT multiplier arguments from Sec. 9.1.

Property (51.44) provides some intuition on how the choice of the penalty factor $q(w)$ defines the solution space. Figure 51.1 plots three contour curves in 2-dimensional space corresponding to the level sets:

$$\|w\|_1 = 1, \quad \|w\|^2 = 1, \quad \|w\|_1 + \|w\|^2 = 1 \quad (51.45)$$

It is seen from the figure that for ℓ_2 -regularization, the search space for w is limited to a region delineated by a circular boundary. In comparison, the search space for ℓ_1 -regularization is delineated by a rotated square boundary with sharp edges, while the search space for elastic-net regularization is midway between these two options. All three regions are obviously convex.

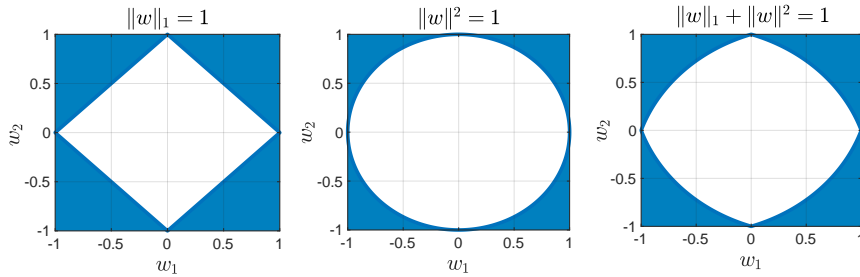


Figure 51.1 The figure illustrates the boundary curves corresponding to conditions (51.45) in 2-dimensional space. The search space for the parameter, w , is limited to the inside of the regions delineated by these curves. Observe that in all three cases, the search domain is convex.

The particular shape for the boundary of the ℓ_1 -region helps promote sparsity, i.e., it helps lead to solutions w_{reg}^* with many zero entries. This is illustrated schematically in Fig. 51.2, which shows boundary curves corresponding to the regions $\|w\|_1 \leq \tau$ and $\|w\|^2 \leq \tau$, along with contour curves for the unregularized risk function, $P(w)$. The solution w_{reg}^* occurs at the location where the contour

curves meet the boundary regions. It is seen, due to the corners that are present in the region $\|w\|_1 \leq \tau$, that the contour curves are more likely to touch this region at a corner point where some of the coordinates are zero. We will establish this conclusion more formally in the next section.

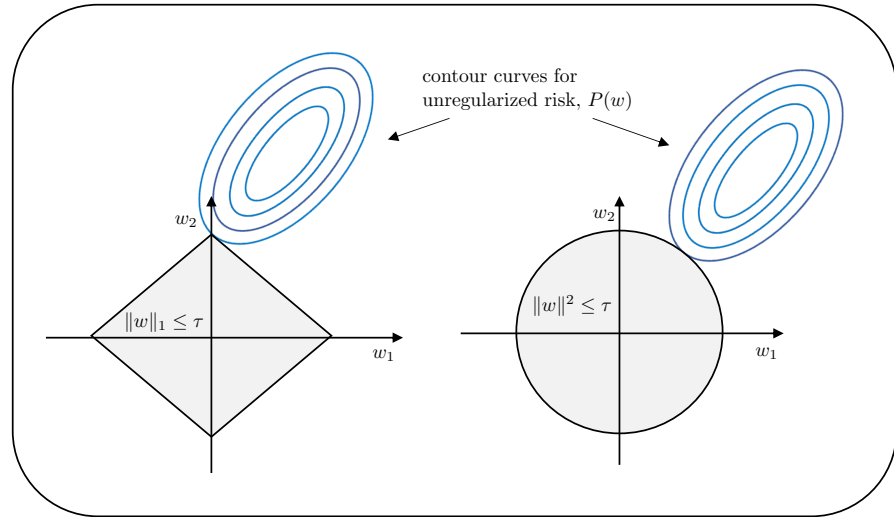


Figure 51.2 Boundary curves corresponding to the regions $\|w\|_1 \leq \tau$ and $\|w\|^2 \leq \tau$, along with contour curves for the unregularized risk function, $P(w)$.

51.4 SOFT THRESHOLDING

We are ready to examine the ability of ℓ_1 -regularization to find sparse solution vectors, w_{reg}^* . A sparse solution helps avoid over-fitting especially for large dimensional data (i.e., when M is large). This is because, when each observation vector y_t has many entries, a sparse w_{reg}^* assigns zero weights to those entries of y_t that are deemed “irrelevant.” For this reason, we say that ℓ_1 -regularization embodies an *automatic* selection capability into the solution by picking only entries from y_t that are most significant to the task of inferring $x(t)$.

For the benefit of the reader, we first review a useful result established earlier in Sec. 11.1.2 and which relies on the soft-thresholding function $\hat{w} = \mathbb{T}_{\frac{\beta}{2}}(z)$. This function operates on the individual entries of its vector argument z to generate the corresponding entries of \hat{w} . For each scalar x , the transformation $\mathbb{T}_{\frac{\beta}{2}}(x)$, with parameter $\beta \geq 0$, is defined as follows:

$$\mathbb{T}_{\frac{\beta}{2}}(x) \triangleq \begin{cases} x - \frac{\beta}{2}, & \text{if } x \geq \frac{\beta}{2} \\ 0, & \text{if } -\frac{\beta}{2} < x < \frac{\beta}{2} \\ x + \frac{\beta}{2}, & \text{if } x \leq -\frac{\beta}{2} \end{cases} \quad (51.46)$$

LEMMA 51.1. (Soft-thresholding operation) Given $z \in \mathbb{R}^M$, a constant $\beta \geq 0$, and a scalar ϕ , the solution to the optimization problem:

$$\hat{w} \triangleq \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \beta \|w\|_1 + \|w - z\|^2 + \phi \right\} \quad (51.47)$$

is unique and given by

$$\hat{w} = \mathbb{T}_{\frac{\beta}{2}}(z) \quad (51.48)$$

The soft-thresholding transformation $\mathbb{T}_{\frac{\beta}{2}}(z)$ helps promote sparse solutions \hat{w} (i.e., solutions with a few nonzero entries). This property is achieved in a measured manner since soft-thresholding sets to zero all entries of z whose magnitude is below the threshold value $\beta/2$, and reduces the size of the larger values by $\beta/2$. Figure 51.3 plots the function $\mathbb{T}_{\frac{\beta}{2}}(x)$ defined by (51.46). In summary, using the ℓ_1 -penalty term in (51.47) results in a sparse solution \hat{w} that is “close” to the vector z .

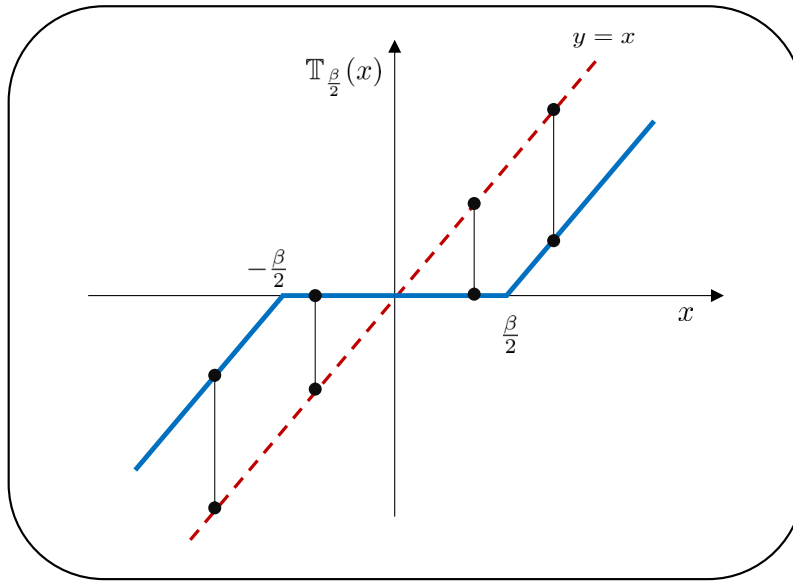


Figure 51.3 The soft-thresholding function, $\mathbb{T}_{\frac{\beta}{2}}(x)$, reduces the value of x gradually. Small values of x within the interval $[-\frac{\beta}{2}, \frac{\beta}{2}]$ are set to zero, while values of x outside this interval have their size reduced by an amount equal to $\beta/2$. The dotted segment represents the line $y = x$.

51.4.1 Orthogonal Data

Before studying the general case of arbitrary data matrices H , we consider first the special case when the “squared matrix” $H^\top H$ happens to be “orthogonal”, namely, when H satisfies

$$H^\top H = \kappa^2 I_M, \quad \text{for some } \kappa^2 > 0 \quad (51.49)$$

Using this normalization condition, and the compact vector and matrix notation $\{d, H\}$ defined in (51.1c), we rewrite the unregularized and regularized risks in the form

$$\begin{aligned} P(w) &\triangleq \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \\ &= \frac{1}{N} \|d - Hw\|^2 \\ &= \frac{1}{N} \left\{ \|d\|^2 - 2d^\top Hw + \kappa^2 \|w\|^2 \right\} \end{aligned} \quad (51.50)$$

and

$$P_{\text{reg}}(w) = \alpha \|w\|_1 + \frac{1}{N} \left\{ \|d\|^2 - 2d^\top Hw + \kappa^2 \|w\|^2 \right\} \quad (51.51)$$

Note that both risks are strongly-convex since $\kappa^2 > 0$. Therefore, they each have a unique global minimizer, denoted by w^\star and w_{reg}^\star .

LEMMA 51.2. (ℓ_1 -regularized solution for orthogonal data) *Consider the ℓ_1 -regularized problem (51.51) under the orthogonality condition (51.49). The solution is unique and given by*

$$w_{\text{reg}}^\star = \mathbb{T}_{\frac{\alpha N}{2\kappa^2}}(w^\star) \quad (51.52)$$

where $w^\star = \frac{1}{\kappa^2} H^\top d$ is the minimizer for the unregularized risk (51.50).

Proof: We employ a completion-of-squares argument to write (51.51) as

$$P_{\text{reg}}(w) = \alpha \|w\|_1 + \frac{\kappa^2}{N} \left\{ \|w\|^2 - \frac{2}{\kappa^2} d^\top Hw + \frac{1}{\kappa^2} \|d\|^2 \right\} \quad (51.53)$$

$$\begin{aligned} &= \alpha \|w\|_1 + \frac{\kappa^2}{N} \left\{ \left\| w - \frac{1}{\kappa^2} H^\top d \right\|^2 + \frac{1}{\kappa^2} \|d\|^2 - \frac{1}{\kappa^4} \|H^\top d\|^2 \right\} \\ &\propto \beta \|w\|_1 + \|w - z\|^2 + \phi \end{aligned} \quad (51.54)$$

where \propto is the proportionality symbol, while the scalars $\{\beta, \phi\}$ and the column vector $z \in \mathbb{R}^M$ are defined by

$$\beta \triangleq \frac{\alpha N}{\kappa^2} > 0 \quad (51.55a)$$

$$z \triangleq \frac{1}{\kappa^2} H^\top d = w^\star \quad (51.55b)$$

$$\phi \triangleq \frac{1}{\kappa^2} \|d\|^2 - \frac{1}{\kappa^4} \|H^\top d\|^2 \quad (51.55c)$$

Observe that z agrees with the minimizer, w^* , for the unregularized problem under condition $H^\top H = \kappa^2 I$. Minimization of the empirical risk (51.54) is now of the same form as problem (51.47). Therefore, we deduce that the minimizer to (51.49) under the orthogonality condition (51.49) is given by (51.52). ■

Observe how construction (51.52) applies soft-thresholding to w^* with the threshold defined by $\alpha N/2\kappa^2$; this value (and, hence, sparsity) increases with α .

51.4.2 LASSO or Basis Pursuit Denoising

More generally, for data matrices H that do not satisfy the orthogonality condition (51.49), we can derive a similar expression for w_{reg}^* involving a soft-thresholding operation, albeit one where w^* is replaced by another vector defined in terms of a *dual* variable — see expression (51.61) further ahead. We will derive the result under both ℓ_1 and elastic-net regularization.

Thus, consider the regularized least-squares problem:

$$w_{\text{reg}}^* = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ P_{\text{reg}}(w) \triangleq q(w) + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\} \quad (51.56a)$$

where the regularization factor has the form

$$q(w) = \alpha \|w\|_1 + \rho \|w\|^2, \quad \alpha > 0, \quad \rho \geq 0 \quad (51.56b)$$

When $\rho = 0$ we have pure ℓ_1 -regularization. Problem (51.56a)–(51.56b) is known as LASSO, where the acronym stands for “*least absolute shrinkage and selection operator*.” The shrinkage feature is because the resulting solution w_{reg}^* will end up satisfying property (51.43). The selection feature is because the same solution will be sparse with generally few nonzero entries. Problem (51.56a) is also known as the *basis pursuit denoising problem*; this is because it seeks a sparse representation for the vector d in terms of the columns of H .

Unfortunately, when H is not orthogonal, a closed-form expression for the solution w_{reg}^* is not possible any longer. For this reason, the LASSO problem (51.56a) is usually solved iteratively by means of subgradient or proximal gradient iterations, with or without stochastic sampling of data, as was already shown earlier in several instances including in Examples 14.1, 15.3, and 16.12; the latter example describes a *stochastic proximal* gradient implementation that relies on instantaneous gradient approximations and which we reproduce here illustration purposes.

Stochastic proximal gradient algorithm for LASSO problem (51.56a)

given dataset $\{x(m), y_m\}_{m=0}^{N-1}$;
 start from an arbitrary initial condition, w_{-1} .
repeat until convergence over $n \geq 0$:
 select at random a sample $(x(n), y_n)$ at iteration n ;
 $z_n = (1 - 2\mu\rho)w_{n-1} + 2\mu y_n(x(n) - y_n^\top w_{n-1})$
 $w_n = \mathbb{T}_{\mu\alpha}(z_n)$
end
 return $w^* \leftarrow w_n$.

(51.57)

Other implementations are of course possible. For instance, Example 15.3 describes a full-batch implementation leading to the iterated soft-thresholding algorithm (ISTA):

$$\begin{cases} z_n = (1 - 2\mu\rho)w_{n-1} + \frac{2\mu}{N} \sum_{m=0}^{N-1} (x(m) - y_m^\top w_{n-1}) \\ w_n = \mathbb{T}_{\mu\alpha}(z_n) \end{cases} \quad (51.58)$$

Numerical solutions of the LASSO optimization problem based on the use of convex optimization packages are also possible. The derivation in this section is meant to highlight some properties of the exact solution, such as showing that it continues to have a soft-thresholding form. To do so, we will follow a duality argument.

Expression for LASSO solution

Using the vector notation $\{d, H\}$, problem (51.56a) can be recast as

$$w_{\text{reg}}^* = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ q(w) + \frac{1}{N} \|d - Hw\|^2 \right\} \quad (51.59)$$

or, equivalently, in terms of an auxiliary variable transformation that introduces a constraint:

$$(w_{\text{reg}}^*, z^*) = \underset{w, z}{\operatorname{argmin}} \left\{ q(w) + \frac{1}{N} \|d - z\|^2 \right\} \quad (51.60)$$

subject to $z = Hw$

where we introduced $z \in \mathbb{R}^N$; it depends linearly on w . The risk function in statement (51.60) is convex over w and z . We therefore have a convex optimization problem with a linear equality constraint. This type of formulation is a special case of problem (9.1), involving convex costs subject to convex inequality and equality constraints, and which we studied in Sec. 9.1. The results from that section show that strong duality holds for problem (51.60). This means that we can

learn about the solution w_{reg}^* by using duality arguments to establish the next theorem for both cases of $\rho \neq 0$ and $\rho = 0$; the proof appears in Appendix 51.B.

THEOREM 51.1. (Expression for LASSO solution) *Consider the regularized problem (51.56a)–(51.56b). The solution is unique and admits the following representation:*

$$w_{\text{reg}}^* = \frac{1}{2\rho} \mathbb{T}_\alpha(H^\top \lambda^\circ) \quad (51.61)$$

where $\lambda^\circ \in \mathbb{R}^N$ is determined as follows:

(a) **(elastic-net regularization, $\rho \neq 0$):** λ° is the unique maximum of the following strongly-concave function:

$$\lambda^\circ = \operatorname{argmax}_{\lambda \in \mathbb{R}^N} \left\{ \lambda^\top d - \frac{N}{4} \|\lambda\|^2 - \frac{1}{4\rho} \|\mathbb{T}_\alpha(H^\top \lambda)\|^2 \right\} \quad (51.62)$$

(b) **(ℓ_1 -regularization, $\rho = 0$):** λ° is the unique projection of the vector $\frac{2}{N}d$ onto the set of vectors λ satisfying $\|H^\top \lambda\|_\infty \leq \alpha$:

$$\lambda^\circ = \operatorname{argmin}_{\lambda \in \mathbb{R}^N} \left\| \lambda - \frac{2}{N}d \right\|^2, \quad \text{subject to } \|H^\top \lambda\|_\infty \leq \alpha \quad (51.63)$$

Comparing expression (51.61) with (51.52) for “orthogonal” data matrices, we note that the argument of the soft-thresholding function is now defined in terms of a dual variable λ° and not in terms of the unregularized solution, w^* . Moreover, the threshold in $\mathbb{T}_\alpha(\cdot)$ increases with α so that more sparse models are expected for larger α . Clearly, solving the LASSO problem via (51.61) is not simpler than solving the original optimization problem (51.56a) because we still need to determine λ° in (51.62) or (51.63). The usefulness of result (51.61) is that it provides a representation for the solution in a manner similar to (51.52) and helps illustrate the sparsity properties of the resulting w_{reg}^* . The parameters α and ρ define the degree of regularization: larger values tend to promote smaller (in norm) and more sparse solutions. One useful way to select these parameters is the cross validation technique described later in Sec. 61.3.

Example 51.5 (Comparing different regularized solutions) In this example we compare numerically the behavior of ℓ_2 , ℓ_1 , and elastic-net regularization solutions. First, however, we need to show how to approximate the regularized solution to (51.56a)–(51.56b). We already know that we can employ a stochastic subgradient algorithm for this purpose to arrive at good approximations for w_{reg}^* . Under elastic-net regularization, the recursion would start from some random initial guess, denoted by w_{-1} , and then iterate as follows:

$$w_n = (1 - 2\mu\rho)w_{n-1} - \mu\alpha \operatorname{sign}(w_{n-1}) + 2\mu y_n(x(n) - y_n^\top w_{n-1}), \quad n \geq 0 \quad (51.64)$$

where μ is a small step-size parameter and the notation w_n denotes the approximation for the regularized solution at iteration n . The sign function, when applied to a vector argument, returns a vector with entries equal to ± 1 depending on the signs of the individual entries of w_{n-1} : $+1$ for nonnegative entries and -1 for negative entries. The

algorithm is run *multiple* times over the training data $\{x(n), y_n\}$, with the data being randomly reshuffled at the beginning of each epoch, namely,

- (a) At the start of each epoch, the data $\{x(n), y_n\}$ is randomly reshuffled so that each epoch runs over the same dataset albeit in a different random order.
- (b) The initial condition for the epoch of index k is the iterate value that was obtained at the end of the previous epoch.

The iterate that is obtained at the end of the last epoch is the one that is taken to be the approximation for w_{reg}^* .

Iteration (51.64) applies to both cases of ℓ_1 -regularization (by setting $\rho = 0$) and elastic-net regularization when both α and ρ are positive. Although we already have a closed-form solution for the ℓ_2 -regularized solution via expression (51.21), or can even arrive at it by means of the recursive least-squares (RLS) algorithm (50.123), the same stochastic recursion (51.64) can be used to approximate the ℓ_2 -regularized solution as well by setting $\alpha = 0$; the recursion leads to a computationally simpler algorithm than RLS albeit at a slower convergence rate.

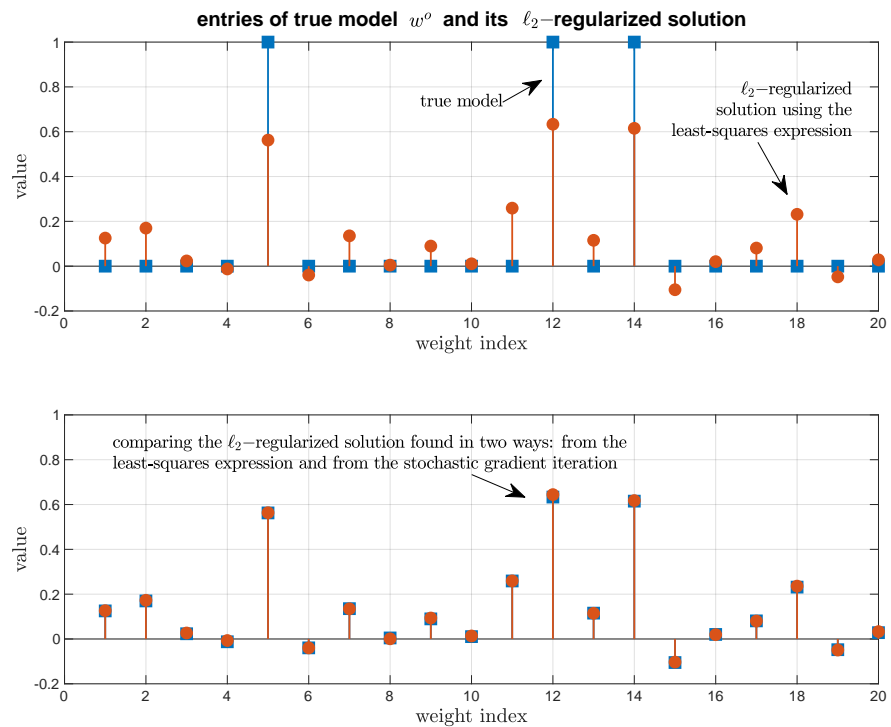


Figure 51.4 The top plot shows the true model w^o with three nonzero entries at value one while all other entries are at zero. The top plot also shows the ℓ_2 -regularized solution, w_{reg}^* , that is obtained by using the least-squares expression (51.21). The bottom plot compares the solutions that are obtained from the least-squares expression (51.21) and from the stochastic recursion (51.64) using 20 runs over the data, $\alpha = 0$, and $\mu = 0.0001$. It is seen that recursion (51.64) is able to learn the ℓ_2 -regularized solution well.

We use the stochastic construction (51.64) to illustrate the behavior of the different

regularization modes, by considering the following numerical example. We generate $N = 4000$ random data points $\{\mathbf{x}(n), \mathbf{y}_n\}$ related through the linear model:

$$\mathbf{x}(n) = \mathbf{y}_n^\top \mathbf{w}^o + \mathbf{v}(n) \quad (51.65)$$

where $\mathbf{v}(n)$ is while Gaussian noise with variance $\sigma_v^2 = 0.01$, and each observation vector has dimension $M = 20$. We generate a sparse true model \mathbf{w}^o consisting of three randomly-chosen entries set to one, while all other entries of \mathbf{w}^o are set to zero. Figures 51.4 and 51.5 illustrate the results that follow from using 20 runs over the data with $\mu = 0.0001$, $\alpha = 5$, and $\rho = 2$. It is seen in the lower plot from the first figure that the stochastic recursion (51.64) converges to a good approximation for the actual least-squares solution from (51.21). The middle plot of the second figure illustrates the sparsity property of the ℓ_1 -regularized solution.

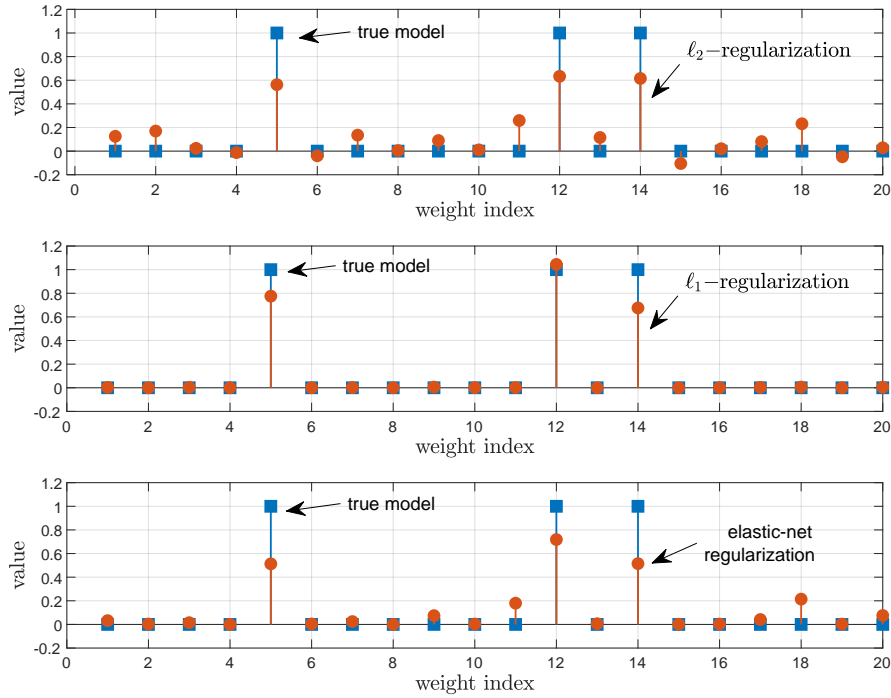


Figure 51.5 All three plots show the true model \mathbf{w}^o with three nonzero entries at value one while all other entries are at zero. In each case, the true model is compared against the ℓ_2 -regularized solution (*top plot*), the ℓ_1 -regularized solution using $\alpha = 5$ (*middle plot*), and the elastic-net regularized solution using $\alpha = 5$ and $\rho = 2$ (*bottom plot*). All these regularized solutions are obtained by using the stochastic (sub)gradient recursion (51.64) using 20 runs over the data and $\mu = 0.0001$. The middle plot illustrates how ℓ_1 -regularization leads to a sparse solution, while the elastic-net regularized solution has slightly more non-zero entries.

51.5 COMMENTARIES AND DISCUSSION

Tikhonov regularization. The regularized least-squares problem (51.16) and its solution (51.21) were proposed by the Russian mathematician **Andrey Tikhonov (1906–1993)** in the publication by Tikhonov (1963) on ill-posed problems — see also the text by Tikhonov and Arsenin (1977). This form of regularization is nowadays very popular and is known as *Tikhonov regularization*. Tikhonov’s formulation was general and applicable to infinite-dimensional operators and not only to finite-dimensional least-squares problems. His work was aimed at solving integral equations of the first-kind, also known as Fredholm integral equations, which deal with the problem of determining a function solution $x(t)$ to an integral equation of the following form:

$$\int_a^b A(s, t)x(t)dt = b(s) \quad (51.66)$$

for a given kernel function, $A(s, t)$, and another function $b(s)$. These integral equations can be ill-conditioned and can admit multiple solutions. The analogy with linear systems of equations of the form $Ax = b$ becomes apparent if we employ the operator notation to rewrite the integral equation in the form $\mathcal{A}x = b$, in terms of some infinite-dimensional operator \mathcal{A} . It turns out that both Phillips (1962) and Tikhonov (1963) proposed using ℓ_2 -regularization to counter ill-conditioning for Fredholm integral equations, which is why this type of regularization is also referred to as the Phillips–Tikhonov or Tikhonov–Phillips regularization. The same technique also appeared in Hoerl (1962), albeit for finite-dimensional operators (i.e., for matrices) in the context of least-squares problems. This latter work was motivated by the earlier contribution on ridge analysis from Hoerl (1959) — see also Hoerl and Kennard (1970) and the review by Hoerl (1985). It is for this reason that ℓ_2 -regularization is also known as ridge regression in the statistics literature. Useful overviews on the role of Tikhonov regularization in the solution of linear systems of equations and least-squares problems appear in the survey article by Neumaier (1998) and in the texts by Golub and Van Loan (1996), Björck (1996), and Hansen (1997). More information on regularization in general can be found in the texts by Wahba (1990) and Engl, Hanke, and Neubauer (1996).

LASSO and basis pursuit denoising. In Examples 51.2 and 51.4 we showed that regularization in the least-squares case corresponds to associating a prior distribution with the sought-after parameter, w (now treated as a random quantity). A Gaussian prior leads to ℓ_2 -regularization, while a Laplacian prior leads to ℓ_1 -regularization as noted by Tibshirani (1996b). We showed in the body of the chapter that ℓ_1 -regularization leads to sparse solutions. However, it has been observed in practice that it tends to retain more nonzero entries than necessary in the solution vector and, moreover, if several entries in the observation space are strongly correlated, the solution vector will tend to keep one of them and discard the others — see Zou and Hastie (2005). Elastic-net regularization, on the other hand, combines ℓ_1 and ℓ_2 -penalty terms and inherits some of their advantages: it promotes sparsity without totally discarding highly correlated observations. This form of regularization was proposed by Zou and Hastie (2005); examination of some of its properties appears in this reference as well as in the text by Hastie, Tibshirani, and Friedman (2009) and in De Mol, De Vito, and Rosasco (2009).

Given data $\{x(n), y_n \in \mathbb{R}^M\}$, the pure ℓ_1 -regularization formulation solves

$$w_{\text{reg}}^* = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ \alpha \|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} \left(x(n) - y_n^\top w \right)^2 \right\} \quad (51.67)$$

where $\alpha > 0$ is the regularization parameter. We explained in the chapter that this

problem is equivalent to solving

$$w_{\text{reg}}^* = \operatorname{argmin}_{w \in \mathbb{R}^M} \|d - Hw\|^2, \quad \text{subject to } \alpha\|w\|_1 \leq \tau \quad (51.68)$$

for some $\tau > 0$. Problems of this type were first proposed by Santosa and Symes (1986) and later by Tibshirani (1996b); the latter reference uses the acronym LASSO for such problems. A similar problem was studied by Chen, Donoho, and Saunders (1998,2001) under the name *basis pursuit denoising*. They examined instead the reverse formulation:

$$w_{\text{reg}}^* = \operatorname{argmin}_{w \in \mathbb{R}^M} \|w\|_1, \quad \text{subject to } \|d - Hw\|^2 \leq \epsilon \quad (51.69)$$

for some small $\epsilon > 0$. This formulation was motivated by the earlier work in Chen and Donoho (1994) on standard *basis pursuit*. In this latter problem, the objective is to seek a sparse representation for a signal vector d from an overcomplete basis H , namely, to solve — see Prob. 51.7:

$$\min_{w \in \mathbb{R}^M} \|w\|_1, \quad \text{subject to } d = Hw \quad (51.70)$$

All three formulations (51.67), (51.68), and (51.69) are equivalent to each other for suitable choices of the parameters $\{\alpha, \tau, \epsilon\}$ — see Prob. 51.7. The contributions by Tibshirani (1996b) and Chen, Donoho, and Saunders (1998,2001) generated renewed interest in ℓ_1 -regularized problems in the statistics, machine learning, and signal processing literature. These types of problems have an older history, especially in the field of geophysics. For example, a problem of the same form as (51.67) was used in the deconvolution of seismic signals by Santosa and Symes (1986). Their work was motivated by the earlier contributions by Claerbout and Muir (1973) and Taylor, Banks, and McCoy (1979). Using our notation, these last two references consider optimization problems of the following form (compare with (51.67)):

$$w_{\text{reg}}^* = \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \alpha\|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} |x(n) - y_n^T w| \right\} \quad (51.71)$$

where the sum of the absolute residuals (rightmost term) is used in place of the sum of their squared values, as is the case in (51.67). Both formulations employ an ℓ_1 -penalty term. One of the earliest recognitions that ℓ_1 -regularization promotes sparsity appears in the article by Santosa and Symes (1986, p. 1308), where it is stated that the use of the ℓ_1 -penalty term “has the effect of constructing a solution which has the least number of nonzero components.” Arguments and derivations in support of the sparsity-promoting property of the ℓ_1 -penalty appear in Levy and Fullagar (1981), Oldenburg, Scheuer, and Levy (1983), and also in Santosa and Symes (1986, Sec. 2). In their formulation of the deconvolution problem, Santosa and Symes (1986) proposed replacing (51.71) by the same problem (51.67) using the sum of squared residuals — see their expressions (1.16) and (5.1).

It is useful to note that design problems involving ℓ_1 -measures of performance have also been pursued in the control field, starting from the mid 1980s. The primary motivation there for the use of the ℓ_1 -norm has been to design control laws that minimize the effect of persistent *bounded* disturbances on the output of the system. Among the earliest references that promoted this approach are the works by Vidyasagar (1986) and Dahleh and Pearson (1986,1987). A thorough treatment of the subject matter, along with an extensive bibliography, appears in the text by Dahleh and Diaz-Bobillo (1995).

Robust least-squares designs. Given an $N \times M$ data matrix H , an $N \times 1$ target vector d , an $N \times N$ positive-definite weighting matrix R , and an $M \times M$ positive-definite regularization matrix Π , the solution to the following regularized weighted least-squares

problem:

$$w^* \triangleq \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ w^\top \Pi w + (d - Hw)^\top R (d - Hw) \right\} \quad (51.72)$$

is unique and given by

$$w^* = (\Pi + H^\top R H)^{-1} H^\top R d \quad (51.73)$$

When the data $\{d, H\}$ are subject to uncertainties, the performance of this solution can deteriorate appreciably. Assume that the actual data matrix that generated the target signal d is $H + \delta H$ and not H , for some small perturbation δH . Then, the above solution w^* , which is designed based on knowledge of the nominal value H , does not take into account the presence of the perturbations in the data. One way to address this problem is to formulate a robust version of the least-squares problem as follows:

$$w^{\text{rob}} \triangleq \operatorname{argmin}_{w \in \mathbb{R}^M} \max_{\{\delta d, \delta H\}} \left\{ w^\top \Pi w + \left((d + \delta d) - (H + \delta H)w \right)^\top R \left((d + \delta d) - (H + \delta H)w \right) \right\} \quad (51.74a)$$

where $\{\delta d, \delta H\}$ represent the unknown perturbations that are assumed to be modeled as follows:

$$\begin{bmatrix} \delta d & \delta H \end{bmatrix} = P \Delta \begin{bmatrix} e_d & E_H \end{bmatrix} \quad (51.74b)$$

where Δ is an arbitrary contraction matrix satisfying $\|\Delta\| \leq 1$ and $\{P, e_d, E_H\}$ are known quantities of appropriate dimensions, e.g., e_d is a column vector. The matrix P is meant to constrain the perturbations to its range space. Problem (51.74a) can be interpreted as a constrained two-game problem, with the designer trying to select an estimate w^{rob} that minimizes the cost while the opponent $\{\delta d, \delta H\}$ tries to maximize the same cost. It turns out that the solution to (51.74a) has the form of a *regularized* least-squares solution albeit one with modified $\{\Pi, R\}$ matrices, as indicated by the following result.

Robust regularized least-squares (Sayed, Nascimento, and Cipparrone (2002)). *Problem (51.74a)–(51.74b) has a unique solution given by*

$$w^{\text{rob}} = \left(\hat{\Pi} + H^\top \hat{R} H \right)^{-1} \left(H^\top \hat{R} d + \hat{\beta} E_H^\top e_d \right) \quad (51.75a)$$

where $\{\hat{\Pi}, \hat{R}\}$ are obtained from $\{\Pi, R\}$ as follows:

$$\hat{\Pi} = \Pi + \hat{\beta} E_H^\top E_H \quad (51.75b)$$

$$\hat{R} = R + R P (\hat{\beta} I_M - P^\top R P)^\dagger P^\top R \quad (51.75c)$$

where the notation \dagger refers to the pseudo inverse of its matrix argument, and the scalar $\hat{\beta}$ is determined by solving

$$\hat{\beta} = \operatorname{argmin}_{\beta \geq \|P^\top R P\|} G(\beta) \quad (51.75d)$$

where the function $G(\beta)$ is defined as follows:

$$G(\beta) \triangleq \|w(\beta)\|_{\hat{\Pi}(\beta)}^2 + \|d - Hw(\beta)\|_{\hat{R}(\beta)}^2 + \beta \|e_d - E_H w(\beta)\|^2 \quad (51.76)$$

where the notation $\|a\|_X^2$ stands for $a^\top X a$ and

$$R(\beta) = R + RP(\beta I - P^\top RP)^\dagger P^\top R \quad (51.77a)$$

$$\Pi(\beta) = \Pi + \beta E_H^\top E_H \quad (51.77b)$$

$$w(\beta) = \left(\hat{\Pi}(\beta) + H^\top R(\beta) H \right)^{-1} \left(H^\top R(\beta) d + \beta E_H^\top e_d \right) \quad (51.77c)$$

We denote the lower bound on β by $\beta_\ell = \|P^\top RP\|$. Compared with the solution (51.73) to the original regularized least-squares problem, we observe that the expression for w^{rob} is distinct in some important ways:

- (a) First, the weighting matrices $\{\Pi, R\}$ are replaced by corrected versions $\{\hat{\Pi}, \hat{R}\}$. These corrections are defined in terms of a scalar $\hat{\beta}$, which is obtained by minimizing $G(\beta)$ over the semi-open interval $[\beta_\ell, \infty)$.
- (b) It was shown by Sayed and Chen (2002) and Sayed, Nascimento, and Cipparrone (2002) that the function $G(\beta)$ has a unique global minimum (and no local minima) over the interval $[\beta_\ell, \infty)$. This means that the determination of $\hat{\beta}$ can be pursued by standard search procedures without worrying about convergence to undesired local minima. Extensive experiments suggest that setting $\hat{\beta} = \lambda \beta_\ell$ (a scaled multiple of the lower bound for some positive λ chosen by the designer) is generally sufficient.
- (c) The right-hand side of (51.75a) contains an additional term $\hat{\beta} E_H^\top e_d$. The expression for w^{rob} can be viewed as the solution to the following extended problem

$$w^{\text{rob}} = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ \begin{bmatrix} 1 & w^\top \end{bmatrix} \begin{bmatrix} \hat{\beta} \|e_d\|^2 & -\hat{\beta} e_d^\top E_H \\ -\hat{\beta} E_H^\top e_d & \hat{\Pi} \end{bmatrix} \begin{bmatrix} 1 \\ w \end{bmatrix} + \|d - Hw\|_{\hat{R}}^2 \right\} \quad (51.78)$$

- (d) For values $\hat{\beta} > \beta_\ell$, the pseudo-inverse operation can be replaced by standard matrix inversion and it holds that

$$\hat{R}^{-1} = R^{-1} - \hat{\beta}^{-1} P P^\top \quad (51.79)$$

Other robust variations of least-squares are possible. For example, model (51.74b) for the perturbations can be replaced by one of the form

$$\|\delta H\| \leq \eta, \quad \|\delta d\| \leq \eta_d \quad (51.80)$$

where the uncertainties are instead assumed to lie within bounded regions determined by the positive scalars $\{\eta, \eta_d\}$. The solution has a similar structure and is described in Chandrasekaran *et al.* (1997,1998) and Sayed, Nascimento, and Cipparrone (2002). A convex optimization approach is described in El Ghaoui and Lebret (1997). Other variations and geometric arguments are described in Sayed, Nascimento, and Chandrasekaran (1998) — see also Probs. 51.19–51.21.

PROBLEMS

51.1 Consider the least-squares problem (51.1a) with a rank deficient H :

$$H = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 0 & 0 \end{bmatrix}, \quad d = \begin{bmatrix} +1 \\ +1 \\ -1 \end{bmatrix}$$

- (a) Verify that all solutions to the normal equations take the form $w^* = \text{col}\{1 - 2b, b\}$ for any $b \in \mathbb{R}$.
- (b) Verify that all vectors in the nullspace of $H^T H$ take the form $p = \text{col}\{-2b, b\}$.
- (c) Verify that the following are two valid solutions:

$$w_1^* = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad w_2^* = \begin{bmatrix} -3 \\ 2 \end{bmatrix}$$

- (d) Consider the test vector $y_t = \text{col}\{2, 2\}$. Compute the $\hat{x}(t)$ that result from both solutions. *Remark.* Observe that the predictions have opposite signs, which is undesirable in applications where the sign of $\hat{x}(t)$ is used to perform classification.
- 51.2** Consider the least-squares problem (51.1a) with an ill-conditioned matrix H :

$$H = \begin{bmatrix} 1 & & \\ & -1 & \\ & & \sqrt{\epsilon} \end{bmatrix}, \quad d = \begin{bmatrix} +1 \\ -1 \\ +1 \end{bmatrix} \quad (51.81)$$

where $\epsilon > 0$ is a small number, and the entries of d are binary variables of the type $x(n) = \pm 1$.

- (a) What is the condition number of $H^T H$?
- (b) Determine the solution w^* to the normal equations.
- (c) Consider the two observation vectors

$$y_1 = \text{col}\{10, 10, 10^{-6}\}, \quad y_2 = \text{col}\{10, 10, -10^{-6}\}$$

where their trailing entries have small size and differ in sign. Predict their target signals $\hat{x}(1)$ and $\hat{x}(2)$. *Remark.* Observe how \hat{x}_2 can become negative for small enough ϵ while \hat{x}_1 is always positive. If the sign of \hat{x} is used to classify the observation vector y , then the vectors $\{y_1, y_2\}$, despite being very close to each other in Euclidean space, will end up being assigned to different classes.

51.3 Let w^* and w_{reg}^* denote solutions to the unregularized and regularized least-squares risks (51.1a) and (51.16), respectively.

- (a) Show that $w_{\text{reg}}^* = (\rho N I + H^T H)^{-1} H^T H w^*$.
- (b) Introduce the eigen-decomposition $H^T H = U \Lambda U^T$, where U is $M \times M$ orthogonal and Λ is diagonal with nonnegative entries $\{\lambda(m)\}$. Let $\bar{w}_{\text{reg}} = U^T w_{\text{reg}}^*$ and $\bar{w} = U^T w^*$ and denote their individual entries by $\{\bar{w}_{\text{reg}}(m), \bar{w}(m)\}$. Verify that

$$\bar{w}_{\text{reg}}(m) = \left(\frac{\lambda(m)}{\rho N + \lambda(m)} \right) \bar{w}(m), \quad m = 1, 2, \dots, M$$

Conclude that $\|w_{\text{reg}}^*\|^2 < \|w^*\|^2$.

51.4 Refer to the minimizers $\{w^*, w_{\text{reg}}^*\}$ for the unregularized and regularized least-squares problems.

- (a) Show that $P_{\text{reg}}(w_{\text{reg}}^*) - P(w^*) = \rho (w^*)^T w_{\text{reg}}^*$.
- (b) Introduce the same transformations $\{\bar{w}, \bar{w}_{\text{reg}}\}$ from Prob. 51.3 and conclude that

$$P_{\text{reg}}(w_{\text{reg}}^*) - P(w^*) = \sum_{m=1}^M \left(\frac{\rho \lambda(m)}{\rho N + \lambda(m)} \right) |\bar{w}(m)|^2$$

- (c) Since generally at least one $\lambda(m) \neq 0$ and $w^* \neq 0$, conclude that $P_{\text{reg}}(w_{\text{reg}}^*) > P(w^*)$.
- (d) Verify that the function $f(\rho) = \rho \lambda / (\rho N + \lambda)$ is non-decreasing in ρ . Conclude that the bias increases with ρ .

51.5 We re-examine the result of Prob. 50.11 for the case of ℓ_2 -regularized least-squares (or ridge regression). Thus, refer again to the stochastic model (50.88) where v has covariance matrix $\sigma_v^2 I_N$ but is not necessarily Gaussian. Introduce the mean-square error risk, $P(w) = \mathbb{E} \|d - Hw\|^2$, where the expectation is over the source of randomness

in **d**. Verify that the ℓ_2 -regularized least-squares solution $\mathbf{w}_{\text{reg}}^*$ given by (51.21) leads to the following average excess risk expression:

$$\mathbb{E} P(\mathbf{w}_{\text{reg}}^*) - P(w^o) = (w^o)^\top H^\top \left(I + \frac{1}{\rho N} H H^\top \right)^{-2} H w^o + \sigma_v^2 \text{Tr} \left[\left(H(\rho N I + H^\top H)^{-1} H^\top \right)^2 \right]$$

Verify that the expression reduces to the result of Prob. 50.11 as $\rho \rightarrow 0$.

51.6 The expression in Prob. 51.5 consists of two terms: the first one depends on $1/\rho$ while the second one varies with ρ . Show that the average excess risk is bounded by

$$\mathbb{E} P(\mathbf{w}_{\text{reg}}^*) - P(w^o) \leq \frac{\rho N}{2} \|w^o\|^2 + \frac{\sigma_v^2}{2\rho N} \text{Tr}(H^\top H)$$

Minimize the bound over ρ and conclude that $\mathbb{E} P(\mathbf{w}_{\text{reg}}^*) - P(w^o) \leq \sigma_v \|w^o\| \sqrt{\text{Tr}(H^\top H)}$. For which value of ρ is this bound attained?

51.7 Consider the ℓ_1 -regularized problem with $\alpha > 0$:

$$\underset{w \in \mathbb{R}^M}{\text{argmin}} \left\{ \alpha \|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\}$$

Using the vector notation (51.1c), show that the problem is equivalent to solving:

$$\underset{w \in \mathbb{R}^M}{\text{argmin}} \|w\|_1, \quad \text{subject to } \|d - Hw\|^2 \leq \epsilon$$

for some $\epsilon \geq 0$. Show that as $\alpha \rightarrow 0$, the formulation reduces to the so-called *basis pursuit* problem (which involves an equality constraint):

$$\underset{w \in \mathbb{R}^M}{\text{argmin}} \|w\|_1, \quad \text{subject to } Hw = d$$

51.8 Establish the validity of expression (51.102) for $\mathbb{S}_\alpha(x)$.

51.9 In this problem, we follow the approach described in the earlier Example 14.10 to express the ℓ_1 -regularized least-squares (LASSO) solution in an alternative form. Consider the regularized problem:

$$w_{\text{reg}}^* = \underset{w \in \mathbb{R}^M}{\text{argmin}} \left\{ P_{\text{reg}}(w) = \alpha \|w\|_1 + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\}$$

We denote the individual entries of w and y_n by $w = \text{col}\{w_m\}$ and $y_n = \text{col}\{y_{n,m}\}$, respectively, for $m = 1, 2, \dots, M$. We also use the notation w_{-m} and $y_{n,-m}$ to refer to the vectors w and y_n with their m -th entries excluded.

(a) Verify that, as a function of w_m , the regularized risk can be written as:

$$P_{\text{reg}}(w) = a_m \left\{ \frac{\alpha}{a_m} |w_m| + \left(w_m - \frac{c_m}{a_m} \right)^2 \right\} + \text{terms indep. of } w_m$$

where

$$a_m \triangleq \frac{1}{N} \sum_{n=0}^{N-1} y_{n,m}^2, \quad c_m \triangleq \frac{1}{N} \sum_{n=0}^{N-1} y_{n,m} (x(n) - y_{n,-m}^\top w_{-m})$$

(b) Conclude that the minimizer over w_m is given by $\hat{w}_m = \mathbb{T}_{\alpha/2a_m}(\hat{c}_m/a_m)$, for $m = 1, 2, \dots, M$, and where \hat{c}_m is given by the same expression as c_m with w_{-m} replaced by \hat{w}_{-m} .

51.10 Replace condition (51.49) by $H^\top H = D^2 > 0$, where D is diagonal. Use result (11.35) to show that expression (51.52) is replaced by

$$w_{\text{reg}}^* = \text{sign}(w^*) \odot \left(|D^2 w^*| - \frac{\alpha N}{2} \right)_+$$

where the operations $\text{sign}(x)$, $|x|$, and $(a)_+$ are applied elementwise.

51.11 Consider the ℓ_2 -regularized risk function $P_{\text{reg}}(w) = \rho \|w\|^2 + P(w)$, where $\rho > 0$ and $P(w)$ is some convex risk in w . Show that $P_{\text{reg}}(w)$ is strongly convex and, therefore, has a unique global minimum.

51.12 Refer to the equivalent problems (51.94).

(a) Assume $q(w) = \rho \|w\|^2$. Show that τ decreases as ρ increases.

(b) Assume $q(w) = \alpha \|w\|_1 + \rho \|w\|^2$, where $\alpha > 0$ and $\rho > 0$. Show that τ decreases as either α or ρ increases.

51.13 Consider the following ℓ_2 -regularized stochastic risk:

$$w_{\text{reg}}^o = \underset{w \in \mathbb{R}^M}{\text{argmin}} \left\{ \rho \|w\|_2 + \mathbb{E}(\mathbf{x} - \mathbf{y}^\top w)^2 \right\}$$

Show that $w_{\text{reg}}^o = R_y(\rho I_M + R_y)^{-1} w^o$, where w^o is the minimizer of the unregularized component, $\mathbb{E}(\mathbf{x} - \mathbf{y}^\top w)^2$.

51.14 Consider the following ℓ_1 -regularized stochastic risk:

$$w_{\text{reg}}^o = \underset{w \in \mathbb{R}^M}{\text{argmin}} \left\{ \alpha \|w\|_1 + \mathbb{E}(\mathbf{x} - \mathbf{y}^\top w)^2 \right\}$$

Assume $R_y = \sigma_y^2 I_M$. Show that $w_{\text{reg}}^o = \mathbb{T}_{\alpha/2\sigma_y^2}(w^o)$, where w^o is the minimizer of the unregularized component, $\mathbb{E}(\mathbf{x} - \mathbf{y}^\top w)^2$.

51.15 Refer to the ℓ_1 -regularized problem (51.38). Verify first that for any scalar $x \in \mathbb{R}$, it holds

$$|x| = \min_{z > 0} \frac{1}{2} \left\{ \frac{x^2}{z} + z \right\}$$

Let w_m denote the individual entries of $w \in \mathbb{R}^M$. Conclude that problem (51.38) can be transformed into

$$\min_{w \in \mathbb{R}^M, \{z_m > 0\}} \left\{ P_{\text{reg}}(w) \triangleq \alpha \sum_{m=1}^M \frac{1}{2} \left(\frac{w_m^2}{z_m} + z_m \right) + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\}$$

Remark. The idea of replacing the regularization factor by smoother forms has been exploited in several works, especially in the context of optimization and image processing — see, for example, Geman Yang (1995), Bach *et al.* (2012), Chan and Liang (2014), and Lanza *et al.* (2015).

51.16 Consider a vector $w \in \mathbb{R}^M$ with individual entries $\{w_m\}$. For any $p \geq 1$ and $\delta \geq 0$, the bridge regression problem in statistics refers to

$$\min_{w \in \mathbb{R}^M} \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2, \quad \text{subject to } \sum_{m=1}^M |w_m|^p \leq \delta$$

Show that this problem is equivalent to solving

$$\min_{w \in \mathbb{R}^M} \left\{ \rho \sum_{m=1}^M |w_m|^p + \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - y_n^\top w)^2 \right\}$$

for some $\rho \geq 0$. That is, show that for any $\delta \geq 0$ there exists a $\rho \geq 0$ that makes both problems equivalent to each other (i.e., have the same solution). *Remark.* See the works by Frank and Friedman (1993) and Fu (1998) for a related discussion.

51.17 Refer to the ℓ_1 -regularized problem (51.38) and define the quantities $\{d, H\}$ shown in (51.1c), where $H \in \mathbb{R}^{N \times M}$. Let $\{u_m\}$ denote the individual columns of H for $m = 1, 2, \dots, M$, where each u_m has size $N \times 1$. Let w_m denote the individual entries of $w \in \mathbb{R}^M$. Show that w^* is a solution of (51.38) if, and only if, for every entry w_m^* it holds that

$$\begin{cases} |u_m^\top(d - Hw^*)| \leq N\alpha/2, & \text{when } w_m^* = 0 \\ u_m^\top(d - Hw^*) = \frac{N\alpha}{2} \text{sign}(w_m^*), & \text{when } w_m^* \neq 0 \end{cases}$$

Remark. See Bach *et al.* (2012) for a related discussion.

51.18 Derive expressions (51.117a)–(51.117b) for the conjugate functions.

51.19 Assume H is full rank and has dimensions $N \times M$ with $N > M$. Consider the regularized least-squares problem:

$$w^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ \rho \|w\|^2 + \|d - Hw\|^2 \right\}, \quad \rho > 0$$

and assume $d \notin \mathcal{R}(H)$ and $H^\top d \neq 0$. Let $\tilde{d} = d - Hw^*$ and introduce the scalar $\eta = \rho \|w^*\| / \|\tilde{d}\|$. Verify that $\eta < \|H^\top d\| / \|d\|$. *Remark.* See Sayed, Nascimento, and Chandrasekaran (1998) for a related discussion.

51.20 The next two problems are extracted from Sayed (2003, 2008). Consider an $N \times M$ full rank matrix H with $N \geq M$, and an $N \times 1$ vector d that does not belong to the column span of H . Let η be a positive real number and consider the set of all matrices δH whose 2-induced norms do not exceed η , $\|\delta H\| \leq \eta$. Now consider the following optimization problem whose solution we denote by w^* :

$$w^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ \max_{\|\delta H\| \leq \eta} \|d - (H + \delta H)w\| \right\}$$

That is, we seek to minimize the maximum residual over the set $\{\|\delta H\| \leq \eta\}$.

- (a) Argue from the conditions of the problem that we must have $N > M$.
- (b) Show that the uncertainty set $\{\|\delta H\| \leq \eta\}$ contains a perturbation δH^o such that d is orthogonal to $(H + \delta H^o)$ if, and only if, $\eta \geq \|H^\top d\| / \|d\|$.
- (c) Show that the above optimization problem has a unique solution at $w^* = 0$ if, and only if, the condition on η in part (b) holds.

Remark. For more details on such robust formulations, see Chandrasekaran *et al.* (1997, 1998), Sayed, Nascimento, and Chandrasekaran (1998), and Sayed, Nascimento, and Cipparrone (2002).

51.21 Consider an $N \times M$ full rank matrix H with $N \geq M$, and an $N \times 1$ vector d that does not belong to the column span of H .

- (a) For any nonzero $M \times 1$ column vector w , show that the following rank-one modification of H continues to have full rank for any positive real number η :

$$H(w) \triangleq H - \eta \frac{d - Hw}{\|d - Hw\|} \frac{w^\top}{\|w\|}$$

- (b) Verify that $\|d - H(w)w\| = \|y - Hw\| + \eta \|w\|$, and that the vectors $d - H(w)w$ and $d - Hw$ are collinear and point in the same direction (that is, one is a positive multiple of the other).
- (c) Show that $\|d - H(w)w\| = \max_{\|\delta H\| \leq \eta} \|d - (H + \delta H)w\|$.
- (d) Show that the optimization problem

$$\min_{w \in \mathbb{R}^M} \max_{\|\delta H\| \leq \eta} \|d - (H + \delta H)w\|$$

has a nonzero solution w^* if, and only if, $\eta < \|H^\top d\| / \|d\|$.

- (e) Show that w^* is a nonzero solution of the optimization problem in part (d) if, and only, if $H^\top(w^*)(d - Hw^*) = 0$. That is, the residual vector $d - Hw^*$ should be orthogonal to the perturbed matrix $H(w^*)$. Show further that this condition is equivalent to $H^\top(w^*)(d - H(w^*)w^*) = 0$.
- (f) Assume two nonzero solutions w_1^* and w_2^* exist that satisfy the orthogonality condition of part (e). Argue that $H^\top(w_2^*)(d - H(w_2^*)w_1^*) = 0$, and conclude that $w_1^* = w_2^*$ so that the solution is unique.

Remark. For further details, see Sayed, Nascimento, and Chandrasekaran (1998).

51.A CONSTRAINED FORMULATIONS FOR REGULARIZATION

In this appendix, we first establish the equivalence between problems (51.16) and (51.26) for ℓ_2 -regularized least-squares, and between problems (51.38)–(51.39) and (51.44) for ℓ_1 and elastic-net regularized least-squares. Then we extend the conclusion to other regularized convex risks, besides least-squares. Although it would have been sufficient to treat the general case right away, we prefer to explain the equivalence in a gradual manner for the benefit of the reader, starting with quadratic risks. To establish the equivalence, we will appeal to the Lagrange and KKT multiplier arguments from Sec. 9.1.

51.A.1 Quadratic Risks

We start with the ℓ_2 -regularized least-squares risk.

ℓ_2 -regularization

To begin with, we identify the smallest value for τ . We already know that the solution to the ℓ_2 -regularized problem (51.16) is given by

$$w_{\text{reg}}^* = (\rho NI_M + H^\top H)^{-1} H^\top d \quad (51.82)$$

Now, consider the constrained problem (51.26) for some $\tau > 0$. The unregularized risk $P(w)$ is quadratic in w and is therefore convex and continuously differentiable. The constraint $\rho\|w\|^2 \leq \tau$ defines a convex set in \mathbb{R}^M . We are therefore faced with the problem of minimizing a convex function over a convex domain. It is straightforward to verify that problems of this type can only have global minima — see the argument after (9.10). For the solution w_{reg}^* defined by (51.82) to be included in the search domain $\rho\|w\|^2 \leq \tau$, it is necessary for the value of τ to satisfy $\tau \geq \rho\|w_{\text{reg}}^*\|^2$. This argument shows that the smallest value for τ is

$$\tau = \rho\|w_{\text{reg}}^*\|^2 = \rho\|(\rho NI_M + H^\top H)^{-1} H^\top d\|^2 \quad (51.83)$$

Actually, the regularized solution, w_{reg}^* , will lie on the boundary of the set $\|w\|^2 \leq \tau/\rho$. At the same time, the constraint set will *exclude* any of the solutions, w^* , to the original unregularized solution from (51.1b). This is because $\|w^*\| > \|w_{\text{reg}}^*\|$, as already revealed by (51.24).

Let w_{cons}^* denote a solution to the constrained problem (51.26) for the above value of τ . We want to verify that this solution agrees with w_{reg}^* . We appeal to the KKT conditions from Sec. 9.1. Note first that problem (51.26) does not involve any equality constraints and has only one inequality constraint of the form

$$g(w) \triangleq \rho\|w\|^2 - \tau \leq 0 \quad (51.84)$$

We introduce the Lagrangian function

$$\mathcal{L}(w, \lambda) = P(w) + \lambda(\rho\|w\|^2 - \tau), \quad \lambda \geq 0 \quad (51.85)$$

and let $w_{\text{reg}}^*(\lambda)$ denote a minimizer for it. Strong duality holds because Slater condition (9.58a) is satisfied, i.e., there exists a \bar{w} such that $g(\bar{w}) < 0$ (e.g., $\bar{w} = 0$). The KKT conditions (9.28a)–(9.28e) then state that $w_{\text{reg}}^*(\lambda)$ agrees with w_{cons}^* if, and only if, the following conditions hold for some scalar λ and $w_{\text{reg}}^*(\lambda)$:

$$\rho\|w_{\text{reg}}^*(\lambda)\|^2 - \tau \leq 0, \quad (\text{feasibility of primal problem}) \quad (51.86a)$$

$$\lambda \geq 0, \quad (\text{feasibility of dual problem}) \quad (51.86b)$$

$$\lambda(\rho\|w_{\text{reg}}^*(\lambda)\|^2 - \tau) = 0, \quad (\text{complementary condition}) \quad (51.86c)$$

$$\nabla_w \left\{ \lambda(\rho\|w\|^2 - \tau) + P(w) \right\}_{w=w_{\text{reg}}^*(\lambda)} = 0 \quad (51.86d)$$

If we select $\lambda = 1$, then $w_{\text{reg}}^*(\lambda) = w_{\text{reg}}^*$ and the KKT conditions are satisfied at these values for $\tau = \rho\|w_{\text{reg}}^*\|^2$. It follows that $w_{\text{cons}}^* = w_{\text{reg}}^*$.

ℓ_1 and elastic-net regularization

We can extend the argument for other regularization factors, such as $q(w) = \alpha\|w\|_1$ or $q(w) = \alpha\|w\|_1 + \rho\|w\|^2$. Let w_{reg}^* denote the minimizer for either regularized risk (51.38) or (51.39); the argument applies to both cases. It follows that the smallest value for τ should be:

$$\tau = q(w_{\text{reg}}^*) \quad (51.87)$$

Let w_{cons}^* denote a solution to the constrained problem (51.44) for the above value of τ . We want to verify that this solution agrees with w_{reg}^* . We again appeal to the KKT conditions from Sec. 9.1. Note first that either problem (51.38) or (51.39) does not involve any equality constraints and has only one inequality constraint of the form

$$g(w) \triangleq q(w) - \tau \leq 0 \quad (51.88)$$

We introduce the Lagrangian function

$$\mathcal{L}(w, \lambda) = P(w) + \lambda(q(w) - \tau), \quad \lambda \geq 0 \quad (51.89)$$

and let $w_{\text{reg}}^*(\lambda)$ denote a minimizer for it. Strong duality holds because Slater condition (9.58a) is satisfied, i.e., there exists a \bar{w} such that $g(\bar{w}) < 0$ (e.g., $\bar{w} = 0$). The KKT conditions (9.28a)–(9.28e) then state that $w_{\text{reg}}^*(\lambda)$ agrees with w_{cons}^* if, and only if, the following conditions hold for some scalar λ and $w_{\text{reg}}^*(\lambda)$:

$$q(w_{\text{reg}}^*(\lambda)) - \tau \leq 0, \quad (\text{feasibility of primal problem}) \quad (51.90a)$$

$$\lambda \geq 0, \quad (\text{feasibility of dual problem}) \quad (51.90b)$$

$$\lambda(q(w_{\text{reg}}^*(\lambda)) - \tau) = 0, \quad (\text{complementary condition}) \quad (51.90c)$$

$$0 \in \partial \left\{ q(w) + P(w) \right\}_{w=w_{\text{reg}}^*(\lambda)} \quad (51.90d)$$

If we select $\lambda = 1$, then $w_{\text{reg}}^*(\lambda) = w_{\text{reg}}^*$ and the KKT conditions are satisfied at these values for $\tau = q(w_{\text{reg}}^*)$. It follows that $w_{\text{cons}}^* = w_{\text{reg}}^*$.

51.A.2 Other Convex Risks

The discussion in the body of the chapter reveals that regularization has several benefits: it resolves ambiguities by ensuring unique solutions and counters ill-conditioning and overfitting. Naturally, these favorable conditions come at the expense of introducing bias: the achievable minimum risk (or training error) is higher under regularization

than it would be in the absence of regularization. These various properties have been established so far for the case of least-squares risks. We argue now that regularization ensures similar properties for other convex risk functions besides quadratic risks. Thus, more generally, we let $P(w)$ denote any convex risk function, differentiable or not, and introduce its regularized version:

$$P_{\text{reg}}(w) \triangleq q(w) + P(w) \quad (51.91)$$

where the penalty term, $q(w)$, is also assumed to be convex in w such as the choices introduced earlier in (51.15). In the chapter, we considered one choice for $P(w)$, namely, the quadratic risk (51.92a). Later, when we study learning algorithms, other convex empirical risks will arise (such as logistic risks, exponential risks, hinge risks and others), in which case the results of the current appendix will be applicable; some of the risks will also be non-differentiable. Examples include empirical risks of the form:

$$P(w) = \frac{1}{N} \sum_{n=0}^{N-1} \left(x(n) - y_n^T w \right)^2, \quad (\text{quadratic risk}) \quad (51.92a)$$

$$P(w) = \frac{1}{N} \sum_{n=0}^{N-1} \ln \left(1 + e^{-x(n)y_n^T w} \right), \quad (\text{logistic risk}) \quad (51.92b)$$

$$P(w) = \frac{1}{N} \sum_{n=0}^{N-1} \max \left\{ 0, -x(n)y_n^T w \right\}, \quad (\text{Perceptron risk}) \quad (51.92c)$$

$$P(w) = \frac{1}{N} \sum_{n=0}^{N-1} \max \left\{ 0, 1 - x(n)y_n^T w \right\}, \quad (\text{hinge risk}) \quad (51.92d)$$

Uniqueness of solution. We focus on general optimization problems of the form (51.91), where $P(w)$ is convex in w (but need not be differentiable) and $q(w)$ is one of the convex penalty terms considered before in (51.15).

The first property to note is that whenever $P(w)$ is convex in w , the ℓ_2 -regularized version (i.e., when $q(w) = \rho \|w\|^2$), will be strongly convex for any $\rho > 0$ and, therefore, $P_{\text{reg}}(w)$ will have a unique global minimum, w_{reg}^* . The strong convexity of $P_{\text{reg}}(w)$ in this case follows from the fact that $\rho \|w\|^2$ is itself strongly-convex — see Prob. 51.11. Therefore, ridge regression ensures a unique global minimizer; a similar conclusion can be established when elastic-net regularization is applied for any convex empirical risk $P(w)$. For ℓ_1 -regularization, a unique global minimizer will be guaranteed when $P(w)$ happens to be strictly or strongly convex; in this case, convexity of $P(w)$ alone is not sufficient because the penalty $\alpha \|w\|_1$ is convex but not strictly convex. In the sequel, we assume that the regularized risk $P_{\text{reg}}(w)$ has a unique global minimizer.

Promoting smaller solutions. Let w^* denote a global minimizer for the unregularized convex risk, $P(w)$. This minimizer need not be unique since $P(w)$ is only assumed to be convex but not necessarily strongly convex. Let w_{reg}^* denote the global minimizer for the regularized risk, $P_{\text{reg}}(w)$. This minimizer is unique. Now, since w_{reg}^* minimizes $P_{\text{reg}}(w)$, we have

$$\begin{aligned} q(w_{\text{reg}}^*) + P(w_{\text{reg}}^*) &\leq q(w^*) + P(w^*) \\ \implies q(w_{\text{reg}}^*) - q(w^*) &\leq P(w^*) - P(w_{\text{reg}}^*) \\ &\stackrel{(a)}{\implies} q(w_{\text{reg}}^*) - q(w^*) \leq 0 \\ \iff q(w_{\text{reg}}^*) &\leq q(w^*) \end{aligned} \quad (51.93)$$

where step (a) is because w^* minimizes the unregularized risk, $P(w)$.

Constrained formulation. We assume the regularized risk has a unique global minimizer.

We also assume that the *Slater condition* (9.58a) holds, i.e., there exists a \bar{w} such that $g(\bar{w}) < 0$, which is equivalent to $q(\bar{w}) < \tau$ where $\tau = q(w_{\text{reg}}^*)$. This can be satisfied, for example, at $\bar{w} = 0$ for the penalty terms considered before in (51.15). Then, the same KKT argument used earlier in this appendix under ℓ_1 and elastic-net regularization shows that the following two problems are equivalent (meaning they have the same solution vectors):

$$\begin{cases} w_{\text{reg}}^* = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \{q(w) + P(w)\} \iff \\ w_{\text{cons}}^* = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} P(w) \text{ subject to } q(w) \leq q(w_{\text{reg}}^*) \end{cases} \quad (51.94)$$

51.B EXPRESSION FOR LASSO SOLUTION

In this appendix, we establish Theorem 51.1 for the solution of the LASSO problem under ℓ_1 and elastic-net regularization using a duality argument patterned after the derivation in Chen, Towfic, and Sayed (2015); other related arguments appear in Mota *et al.* (2012,2013).

We assume first that $\rho \neq 0$. We start by introducing the Lagrangian function:

$$\mathcal{L}(w, z, \lambda) \triangleq \frac{1}{N} \|d - z\|^2 + q(w) + \lambda^\top (z - Hw) \quad (51.95)$$

where $\lambda \in \mathbb{R}^N$ is the dual variable (or Lagrange multiplier). The dual function is defined by minimizing $\mathcal{L}(w, z, \lambda)$ over $\{w, z\}$:

$$\begin{aligned} \mathcal{D}(\lambda) &\triangleq \min_{w, z} \mathcal{L}(w, z, \lambda), \quad (\text{dual function}) \\ &= \min_z \left\{ \frac{1}{N} \|d - z\|^2 + \lambda^\top z \right\} + \min_w \left\{ q(w) - \lambda^\top Hw \right\} \end{aligned} \quad (51.96)$$

where we are grouping separately the terms that depend on z and w . Once this dual function is determined, as shown by future expression (51.107), maximizing it leads to the optimal value for λ — see Eq. (51.62):

$$\lambda^o = \underset{\lambda \in \mathbb{R}^N}{\operatorname{argmax}} \mathcal{D}(\lambda) \quad (51.97)$$

Strong duality will then imply that we can determine the optimal solutions for $\{w, z\}$ for problem (51.60) by using this value for λ^o , namely, by solving:

$$z^o \triangleq \underset{z \in \mathbb{R}^N}{\operatorname{argmin}} \left\{ \frac{1}{N} \|d - z\|^2 + (\lambda^o)^\top z \right\} \implies z^o = d - \frac{N}{2} \lambda^o \quad (51.98a)$$

$$w_{\text{reg}}^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ q(w) - (\lambda^o)^\top Hw \right\} \quad (51.98b)$$

Expression (51.98a) shows how z^o is determined from λ^o . We still need to show how to solve (51.98b) and determine the regularized solution in terms of λ^o . We can pursue

this task by appealing to result (51.47). Indeed, note that

$$\begin{aligned}
 w_{\text{reg}}^* &= \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ q(w) - (\lambda^o)^\top H w \right\} \\
 &= \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \alpha \|w\|_1 + \rho \|w\|^2 - (\lambda^o)^\top H w \right\} \\
 &= \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \frac{\alpha}{\rho} \|w\|_1 + \|w\|^2 - \frac{1}{\rho} (\lambda^o)^\top H w \right\} \\
 &= \operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ \frac{\alpha}{\rho} \|w\|_1 + \left\| w - \frac{1}{2\rho} H^\top \lambda^o \right\|^2 - \frac{1}{(2\rho)^2} \|H^\top \lambda^o\|^2 \right\} \quad (51.99)
 \end{aligned}$$

Using result (51.47) we conclude that (51.61) holds.

Determining the dual variable λ^o . To complete the argument, we still need to determine λ^o , which is the maximizer for the dual function $\mathcal{D}(\lambda)$ defined by (51.96). We first determine $\mathcal{D}(\lambda)$. From (51.96) we observe that we need to minimize two separate terms, one over z and one over w .

We already know from the above argument that for any λ :

$$\operatorname{argmin}_{w \in \mathbb{R}^M} \left\{ q(w) - \lambda^\top H w \right\} \implies w_\lambda^* = \frac{1}{2\rho} \mathbb{T}_\alpha(H^\top \lambda) \quad (51.100)$$

where we are denoting the minimizer for a generic λ by the notation w_λ^* . Consequently, the minimum value for this first minimization is given by

$$\begin{aligned}
 q(w_\lambda^*) - \lambda^\top H w_\lambda^* &= \alpha \|w_\lambda^*\|_1 + \rho \|w_\lambda^*\|^2 - \lambda^\top H w_\lambda^* \\
 &\stackrel{(51.100)}{=} \frac{1}{2\rho} \left(\alpha \|\mathbb{T}_\alpha(H^\top \lambda)\|_1 + \frac{1}{2} \|\mathbb{T}_\alpha(H^\top \lambda)\|^2 - \lambda^\top H \mathbb{T}_\alpha(H^\top \lambda) \right)
 \end{aligned} \quad (51.101)$$

To simplify the notation, we let for any vector x :

$$\mathbb{S}_\alpha(x) \triangleq -\alpha \|\mathbb{T}_\alpha(x)\|_1 - \frac{1}{2} \|\mathbb{T}_\alpha(x)\|^2 + x^\top \mathbb{T}_\alpha(x) \quad (51.102)$$

Then, it is verified in Prob. 51.8 that

$$\mathbb{S}_\alpha(x) = \frac{1}{2} \|\mathbb{T}_\alpha(x)\|^2 \quad (51.103)$$

In this way, we can rewrite the minimum value (51.101) more compactly as

$$q(w_\lambda^*) - \lambda^\top H w_\lambda^* = -\frac{1}{2\rho} \mathbb{S}_\alpha(H^\top \lambda) = -\frac{1}{4\rho} \|\mathbb{T}_\alpha(H^\top \lambda)\|^2 \quad (51.104)$$

For illustration purposes, Figure 51.6 plots the soft-thresholding functions $\mathbb{T}_\alpha(x)$ and $\mathbb{S}_\alpha(x)$ for $\alpha = 1$ and a scalar argument, x .

Let us now consider the first minimization in (51.96) for any λ :

$$\min_{z \in \mathbb{R}^N} \left\{ \frac{1}{N} \|d - z\|^2 + \lambda^\top z \right\} \implies \hat{z}_\lambda = d - \frac{N}{2} \lambda \quad (51.105)$$

where we are denoting the minimizer for a generic λ by the notation \hat{z}_λ . Consequently, the minimum value for this minimization is given by

$$\frac{1}{N} \|d - \hat{z}_\lambda\|^2 + \lambda^\top \hat{z}_\lambda = \lambda^\top d - \frac{N}{4} \|\lambda\|^2 \quad (51.106)$$

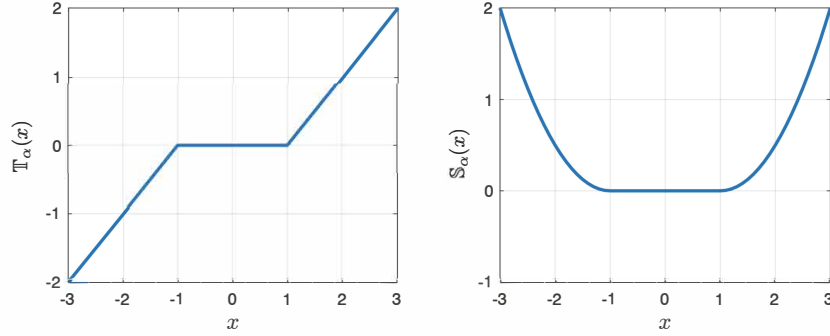


Figure 51.6 Plots of the soft-thresholding functions $\mathbb{T}_\alpha(x)$ and $\mathbb{S}_\alpha(x)$ for $\alpha = 1$.

Adding this result to (51.104) we find that the dual function is given by

$$\mathcal{D}(\lambda) = \lambda^\top d - \frac{N}{4} \|\lambda\|^2 - \frac{1}{4\rho} \left\| \mathbb{T}_\alpha \left(H^\top \lambda \right) \right\|^2 \quad (51.107)$$

It can be verified that this function is strongly-concave and, therefore, has a unique maximum (see next example). The desired dual variable, λ° , is therefore given by (51.62).

The proof technique used so far requires $\rho > 0$. This condition was used to complete the squares in step (51.99). We now explain how to handle the situation $\rho = 0$, which corresponds to pure ℓ_1 -regularization. The solution will continue to be given by expression (51.61), except that λ° will be found by solving the projection problem (51.112). The details are as follows.

We revisit step (51.99) when $\rho = 0$ and note that it reduces to solving a problem of the form:

$$\min_{w \in \mathbb{R}^M} \left\{ \alpha \|w\|_1 - \lambda^\top H w \right\} \quad (51.108)$$

Let \mathcal{C} denote the convex set of vectors satisfying $\|x\|_\infty \leq 1$. We established earlier in Table 8.4 and Prob. 8.55 the following conjugate pair:

$$r(w) = \|w\|_1 \implies r^*(x) = \mathbb{I}_{\mathcal{C}, \infty}[x] \quad (51.109)$$

where the notation $\mathbb{I}_{\mathcal{C}, \infty}[x]$ represents the indicator function relative to set \mathcal{C} : it assumes the value zero if $x \in \mathcal{C}$ and $+\infty$ otherwise. In light of definition (51.114) for the conjugate function, we find that the minimum value of problem (51.108) is given by

$$\min_{w \in \mathbb{R}^M} \left\{ \alpha \|w\|_1 - \lambda^\top H w \right\} = -\mathbb{I}_{\mathcal{C}, \infty}[H^\top \lambda / \alpha] \quad (51.110)$$

Adding this value to (51.106) we find that the dual function is now given by

$$\mathcal{D}(\lambda) = \lambda^\top d - \frac{N}{4} \|\lambda\|^2 - \mathbb{I}_{\mathcal{C}, \infty}[H^\top \lambda / \alpha] \quad (51.111)$$

Maximizing $\mathcal{D}(\lambda)$ over λ results in λ° . To do so, we complete the squares over λ to find that the maximization of $\mathcal{D}(\lambda)$ is equivalent to solving:

$$\lambda^\circ = \operatorname{argmin}_{\lambda \in \mathbb{R}^N} \left\| \lambda - \frac{2}{N} d \right\|^2, \quad \text{subject to } \|H^\top \lambda\|_\infty \leq \alpha \quad (51.112)$$

The minimizer λ^o is obtained by projecting $\frac{2}{N}d$ onto the set of all vectors λ satisfying $\|H^\top \lambda\|_\infty \leq \alpha$. Using $z = Hw$ and $z^o = d - \frac{N}{2}\lambda^o$, we conclude that the optimal solution w_{reg}^* also satisfies the equation

$$Hw_{\text{reg}}^* = d - \frac{N}{2}\lambda^o = \frac{2}{N} \underbrace{\left(\frac{2}{N}d - \lambda^o\right)}_{\text{residual}} \quad (51.113)$$

in terms of the residual resulting from projecting $\frac{2}{N}d$ onto the set $\|H^\top \lambda\|_\infty \leq \alpha$.

Example 51.6 (Duality and conjugate functions) There is an alternative way to arrive at the same expression (51.62) by calling upon the concept of *conjugate* functions, also called Fenchel conjugate functions. This alternative argument is useful for situations (other than least-squares) when explicit expressions for the individual minimum values (51.104) and (51.106) may not be directly available but can be expressed in terms of conjugate functions.

We first recall the definition of conjugate functions from (8.83). Consider an arbitrary function $r(w) : \mathbb{R}^M \rightarrow \mathbb{R}$ with domain $\text{dom}(r)$; the function $r(w)$ need not be convex. Its conjugate function is denoted by $r^*(\lambda) : \mathbb{R}^M \rightarrow \mathbb{R}$ and is defined as:

$$r^*(\lambda) \triangleq \sup_{w \in \mathbb{R}^M} \left\{ \lambda^\top w - r(w) \right\}, \quad \lambda \in \mathcal{Y} \quad (51.114)$$

where \mathcal{Y} denotes the set of all λ where the supremum operation is finite. It can be verified that $r^*(\lambda)$ is convex regardless of whether $r(w)$ is convex or not. Likewise, the set \mathcal{Y} is a convex set — recall Prob. 8.47 and Table 8.4. If $r(w)$ happens to be strongly-convex, then $\mathcal{Y} = \mathbb{R}^M$ (i.e., the sup is finite for all λ).

Now, consider the quadratic function $f(w) = \|w\|^2$ and observe that the dual function $\mathcal{D}(\lambda)$ in (51.96) can be written as:

$$\begin{aligned} \mathcal{D}(\lambda) &= - \sup_{z \in \mathbb{R}^N} \left(-\lambda^\top z - \frac{1}{N} \|d - z\|^2 \right) - \sup_{w \in \mathbb{R}^M} \left(\lambda^\top Hw - q(w) \right) \\ &= - \sup_{z \in \mathbb{R}^N} \left(\lambda^\top (d - z) - \frac{1}{N} \|d - z\|^2 - \lambda^\top d \right) - q^*(H^\top \lambda) \\ &= - \sup_{s \in \mathbb{R}^N} \left(\lambda^\top s - \frac{1}{N} \|s\|^2 - \lambda^\top d \right) - q^*(H^\top \lambda), \quad s \triangleq d - z \\ &= - \sup_{s \in \mathbb{R}^N} \left(\lambda^\top s - \frac{1}{N} \|s\|^2 \right) + \lambda^\top d - q^*(H^\top \lambda) \\ &= - \frac{1}{N} \sup_{s \in \mathbb{R}^N} \left(N\lambda^\top s - \|s\|^2 \right) + \lambda^\top d - q^*(H^\top \lambda) \\ &= - \frac{1}{N} f^*(N\lambda) + \lambda^\top d - q^*(H^\top \lambda) \end{aligned} \quad (51.115)$$

where $f^*(\lambda)$ and $q^*(\lambda)$ denote the conjugate functions of

$$f(w) = \|w\|^2, \quad q(w) = \alpha \|w\|_1 + \rho \|w\|^2 \quad (51.116)$$

Both functions, $f(w)$ and $q(w)$, are strongly-convex and, therefore, the domains of their conjugate functions are the entire space, \mathbb{R}^M . Moreover, since $f(w)$ and $q(w)$ are strongly-convex and differentiable, it follows from the properties of conjugate functions that $f^*(\lambda)$ and $q^*(\lambda)$ are themselves strongly-convex and differentiable — recall Table 8.4. This implies that $\mathcal{D}(\lambda)$ is strongly-concave (i.e., its negative is strongly-convex),

differentiable, and has a unique maximizer, λ° .

It can be verified that the conjugate functions for $f(w)$ and $q(w)$ are given by — see Prob. 51.18:

$$f(w) = \|w\|^2 \implies f^*(\lambda) = \frac{1}{4}\|\lambda\|^2 \quad (51.117a)$$

$$q(w) = \alpha\|w\|_1 + \rho\|w\|^2 \implies q^*(\lambda) = \frac{1}{4\rho}\|\mathbb{T}_\alpha(\lambda)\|^2 \quad (51.117b)$$

Substituting into (51.115), we find that the dual function is given by (51.107).

REFERENCES

- Bach, F., R. Jenatton, J. Mairal, and G. Obozinski (2012), “Optimization with sparsity-inducing penalties,” *Foundations and Trends in Machine Learning*, vol. 4, no. 1, pp. 1–106.
- Bjorck, A. (1996), *Numerical Methods for Least Squares Problems*, SIAM, PA.
- Chan, R. H. and H. X. Liang (2014), “Half-quadratic algorithm for $\ell_p - \ell_q$ problems with applications to TV-1 image restoration and compressive sensing,” pp. 78–103, *Lecture Notes in Comput. Sci.* 8293, Springer, Berlin.
- Chandrasekaran, S., G. Golub, M. Gu, and A. H. Sayed (1997), “Parameter estimation in the presence of bounded modeling errors,” *IEEE Signal Processing Letters*, vol. 4, no. 7, pp. 195–197.
- Chandrasekaran, S., G. Golub, M. Gu, and A. H. Sayed (1998), “Parameter estimation in the presence of bounded data uncertainties,” *SIAM. J. Matrix Anal. Appl.*, vol. 19, no. 1, pp. 235–252.
- Chen, S. and D. Donoho (1994), “Basis pursuit,” *Proc. Asilomar Conference on Signals, Systems and Computers*, pp. 41–44, Pacific Grove, CA.
- Chen, S. S., D. L. Donoho, and M. A. Saunders (1998), “Atomic decomposition by basis pursuit,” *SIAM Journal on Scientific Computing*, vol. 20, no. 1 pp. 33–61. Republished in *SIAM Review*, vol. 43, no. 1, pp. 129–159, 2001. An earlier draft has been available earlier since 1995 as a technical report, Department of Statistics, Stanford University.
- Chen, S., D. Donoho, and M. Saunders (2001), “Atomic decomposition by basis pursuit,” *SIAM Review*, vol. 43, no. 1, pp. 129–159.
- Chen, J., Z. J. Towfic, and A. H. Sayed (2015), “Dictionary learning over distributed models,” *IEEE Trans. Signal Process.*, vol. 63, no. 4, pp. 1001–1016.
- Claerbout, J. F. and F. Muir (1973), “Robust modeling with erratic data,” *Geophysics*, vol. 38, no. 5, pp. 826–844.
- Dahleh, M. A. and I. Diaz-Bobillo (1995), *Control of Uncertain Systems: A Linear Programming Approach*, Prentice Hall, NY.
- Dahleh, M. A. and J. B. Pearson (1986), “ ℓ_1 –optimal feedback controllers for discrete-time systems,” *Proc. Amer. Control Conf. (ACC)*, pp. 1964–1968, Seattle, WA.
- Dahleh, M. A. and J. B. Pearson (1987), “ ℓ_1 –optimal feedback controllers for MIMO discrete-time systems,” *IEEE Trans. Automat. Contr.*, vol. 32, pp. 314–322.
- De Mol, C., E. De Vito, and L. Rosasco (2009), “Elastic-net regularization in learning theory,” *J. Complexity*, vol. 25, no. 2, pp. 201–230.
- El Ghaoui, L. and H. Lebrete (1997), “Robust solutions to least-squares problems with uncertain data,” *SIAM. J. Matrix Anal. Appl.*, vol. 18, no. 4, pp. 1035–1064.
- Engl, H. W., M. Hanke, and A. Neubauer (1996), *Regularization of Inverse Problems*, Kluwer.

- Frank, I. E. and J. H. Friedman (1993), "A statistical view of some chemometrics regression tools," *Technometrics*, vol. 35, pp. 109–148.
- Fu, W. J. (1998), "Penalized regressions: The bridge versus the Lasso," *J. Computational and Graphical Statistics*, vol. 7, no. 3, pp. 397–416.
- Geman, D. and C. Yang (1995), "Nonlinear image recovery with half-quadratic regularization," *IEEE Trans. Image Processing*, vol. 4, no. 7, pp. 932–946.
- Golub, G. H. and C. F. Van Loan (1996), *Matrix Computations*, 3rd edition, The John Hopkins University Press, MD.
- Hansen, P. C. (1997), *Rank-Deficient and Discrete Ill-Posed Problems*, SIAM, PA.
- Hastie, T., R. Tibshirani, and J. Friedman (2009), *The Elements of Statistical Learning*, 2nd edition, Springer, NY.
- Hoerl, A. E. (1959), "Optimum solution of many variables equations," *Chemical Engineering Progress*, vol. 55, pp. 69–78.
- Hoerl, A. E. (1962), "Application of ridge analysis to regression problems," *Chemical Engineering Progress*, vol. 58, pp. 54–59.
- Hoerl, A. E. and R. W. Kennard (1970), "Ridge regression: Biased estimation for nonorthogonal problems," *Technometrics*, vol. 12, no. 1, pp. 55–67.
- Hoerl, R. W. (1985), "Ridge analysis 25 years later," *The American Statistician*, vol. 39, pp. 186–192.
- Lanza, A., S. Morigi, L. Reichel, and F. Sgallari (2015), "A generalized Krylov subspace method for $\ell_p - \ell_q$ minimization," *SIAM J. Sci. Computing*, vol. 37, no. 5, pp. 30–50.
- Levy, S. and P. K. Fullagar (1981), "Reconstruction of a sparse spike train from a portion of its spectrum and application to high-resolution deconvolution," *Geophysics*, vol. 46, no. 9, pp. 1235–1243.
- Mota, J., J. Xavier, P. Aguiar, and M. Puschel (2012), "Distributed basis pursuit," *IEEE Trans. Signal Process.*, vol. 60, no. 4, pp. 1942–1956, Apr. 2012.
- Mota, J., J. Xavier, P. Aguiar, and M. Puschel (2013), "D-ADMM: A communication-efficient distributed algorithm for separable optimization," *IEEE Trans. Signal Process.*, vol. 61, no. 10, pp. 2718–2723.
- Neumaier, A. (1998), "Solving ill-conditioned and singular linear systems: A tutorial on regularization," *SIAM Review*, vol. 40, no. 3, pp. 636–666.
- Oldenburg, D., W. T. Scheuer, and S. Levy (1983), "Recovery of the acoustic impedance from reflection seismograms," *Geophysics*, vol. 48, no. 10, pp. 1318–1337.
- Phillips, D. L. (1962), "A technique for the numerical solution of certain integral equations of the first kind," *Journal of the ACM*, vol. 9, no. 1, pp. 84–97.
- Santosa, F. and W. W. Symes (1986), "Linear inversion of band-limited reflection seismograms," *SIAM J. Sci. Statist. Comput.*, vol. 7, no. 4, pp. 1307–1330.
- Sayed, A. H. and H. Chen (2002), "A uniqueness result concerning a robust regularized least-squares solution," *Systems and Control Letters*, vol. 46, pp. 361–369.
- Sayed, A. H., V. Nascimento, and F. A. M. Cipparrone (2002), "A regularized robust design criterion for uncertain data," *SIAM J. Matrix Anal. Appl.*, vol. 23, no. 4, pp. 1120–1142.
- Sayed, A. H., V. H. Nascimento, and S. Chandrasekaran (1998), "Estimation and control with bounded data uncertainties," *Linear Algebra and Its Applications*, vol. 284, pp. 259–306.
- Taylor, H. L., S. C. Banks, and J. F. McCoy (1979), "Deconvolution with the ℓ_1 norm," *Geophysics*, vol. 44, no. 1, pp. 39–52.
- Tibshirani, R. (1996b), "Regression shrinkage and selection via the Lasso," *J. Royal Stat. Society, Series B*, vol. 58, no. 1, pp. 267–288.
- Tikhonov, A. N. (1963), "Solution of incorrectly formulated problems and the regularization method," *Soviet Math. Dokl.*, vol. 4, pp. 1035–1038.
- Tikhonov, A. N. and V. Y. Arsenin (1977), *Solutions of Ill-Posed Problems*, Winston, NY.
- Vidyasagar, M. (1986), "Optimal rejection of persistent bounded disturbances," *IEEE Trans. Aut. Control*, vol. 31, no. 6, pp. 527–534.
- Wahba, G. (1990), *Spline Models for Observational Data*, SIAM, PA.

Zou, H. and T. Hastie (2005), “Regularization and variable selection via the elastic net,” *J. Royal Stat. Society*, Series B, pp. 301–320.