

Regression and Regularization

For the regression problem, we observe $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$, and try to estimate a function $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$. We fit this function by trading off two factors:

1. **Data Fidelity.** Our solution should satisfy

$$f(\mathbf{x}_i) \approx y_i, \quad \text{for } i = 1, \dots, n$$

This is typically quantified using a **loss function**, which penalizes the deviations of $f(\mathbf{x}_i)$ from y_i . This typically has the form of a single scalar function that is applied to each data point and then added up:

$$\text{Loss}(f, \{\mathbf{x}_i\}_{i=1}^n, \{y_i\}_{i=1}^n) = \sum_{i=1}^n L(f(\mathbf{x}_i), y_i)$$

One commonly used loss function is *squared-error*,

$$L(f(\mathbf{x}_i), y_i) = (y_i - f(\mathbf{x}_i))^2.$$

2. **Modeling and Regularization.** We can temper the regression function in one of two ways. The first is to simply restrict it to lying in some function class \mathcal{F} . We then solve

$$\underset{f \in \mathcal{F}}{\text{minimize}} \sum_{i=1}^n L(f(\mathbf{x}_i), y_i).$$

For example, we might consider the set of *linear functions*:

$$\mathcal{F} = \{f : f(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{x} + \beta_0 \text{ for some } \boldsymbol{\beta} \in \mathbb{R}^d, \beta_0 \in \mathbb{R}\}.$$

The problem of estimating the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is now distilled down to the problem of estimating a vector $\boldsymbol{\theta} = \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix} \in \mathbb{R}^{d+1}$.

There are of course trade-offs in choosing \mathcal{F} . A larger \mathcal{F} gives us a richer class of functions to choose from, but we run the risk of overfitting — the empirical risk $\frac{1}{n} \sum_{i=1}^n L(f^*(\mathbf{x}_i), y_i)$ might be very different than the true risk $\mathbb{E}[L(f^*(X), Y)]$ of our choice f^* .

The second way to mitigate the danger of overfitting is to use a large, rich set \mathcal{F} , but then have some penalty on the complexity of the choices of f inside of this class. This “complexity” is quantified using a regularization function $r(f)$ — there are many choices of r we might consider. We then solve

$$\underset{f \in \mathcal{F}}{\text{minimize}} \sum_{i=1}^n L(f(\mathbf{x}_i), y_i) + \lambda r(f),$$

where $\lambda \geq 0$ is a user-specified parameter that controls the balance between these two terms.

Example: Linear regression using the LASSO

Let's look carefully at one particular (and popular!) technique for regression. We will use the squared-error loss along with \mathcal{F} as the set of linear (plus offset) functions on \mathbb{R}^d :

$$L(f(\mathbf{x}_i), y_i) = (y_i - \boldsymbol{\beta}^T \mathbf{x}_i - \beta_0)^2 = (y_i - \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i)^2,$$

where

$$\boldsymbol{\theta} = \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix}, \quad \tilde{\mathbf{x}}_i = \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix}.$$

Stacking the $\tilde{\mathbf{x}}_i^T$ up as rows of a $n \times (d + 1)$ matrix \mathbf{A} and collecting the y_i into a single vector,

$$\mathbf{A} = \begin{bmatrix} 1 & \mathbf{x}_1^T \\ 1 & \mathbf{x}_2^T \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

we can write

$$\sum_{i=1}^n L(f(\mathbf{x}_i), y_i) = \|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2.$$

For the penalty, we use the sum of the absolute values of the entries in $\boldsymbol{\theta}$:

$$r(f) = \sum_{i=1}^{d+1} |\theta(i)| = \|\boldsymbol{\theta}\|_1.$$

Our optimization program is now¹

$$\text{(LASSO)} \quad \underset{\boldsymbol{\theta}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1.$$

¹We have introduced the factor of 1/2 in front of the loss simply for convenience and consistency with the literature.

This regression technique is commonly referred to as the **LASSO**.

The motivation for using the ℓ_1 norm as a regularizer is that doing so tends to produce $\boldsymbol{\theta}$ which have a small number of non-zero terms. This is certainly true in practice, and we can do a little bit of analysis that explains why.

Specifically, we can show that there is a solution to (LASSO) above that has at most n non-zero entries using a relatively simple argument. Let $\boldsymbol{\theta}$ be any vector with more than n non-zero terms in it:

$$\text{nnz}(\boldsymbol{\theta}) := \#\{i : \theta(i) \neq 0\} \geq n + 1.$$

Then at least one of the following is true:

1. There is another $\boldsymbol{\theta}' \in \mathbb{R}^{d+1}$ such that $\text{nnz}(\boldsymbol{\theta}') \leq \text{nnz}(\boldsymbol{\theta})$ and

$$\frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}'\|_2^2 + \lambda\|\boldsymbol{\theta}'\|_1 < \frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda\|\boldsymbol{\theta}\|_1,$$

or

2. there is another $\boldsymbol{\theta}' \in \mathbb{R}^{d+1}$ such that $\text{nnz}(\boldsymbol{\theta}') < \text{nnz}(\boldsymbol{\theta})$ and

$$\frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}'\|_2^2 + \lambda\|\boldsymbol{\theta}'\|_1 = \frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda\|\boldsymbol{\theta}\|_1.$$

Given $\boldsymbol{\theta}$ as above, let Γ be the locations of the non-zero terms in $\boldsymbol{\theta}$:

$$\Gamma = \{i : \theta(i) \neq 0\}.$$

We are supposing that $|\Gamma| \geq n + 1$. Since $|\Gamma|$ is greater than the number of rows in \mathbf{A} , there is at least one vector \mathbf{z} that is also supported on Γ ,

$$z(i) = 0, \quad \text{for } i \notin \Gamma,$$

such that $\alpha \mathbf{z} \in \text{Null}(\mathbf{A})$ for all $\alpha \in \mathbb{R}$. We will show that by adding a little bit of \mathbf{z} to $\boldsymbol{\theta}$, we can hold the ℓ_2 loss term constant while either decreasing the ℓ_1 regularization term or driving one of the non-zero terms to zero. Notice that

$$\|\mathbf{y} - \mathbf{A}(\boldsymbol{\theta} + \epsilon \mathbf{z})\|_2^2 = \|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2$$

for all $\epsilon > 0$, since $\mathbf{z} \in \text{Null}(\mathbf{A})$.

Now let \mathbf{s}_θ be the vector supported on Γ that contains the signs of the non-zero entries of $\boldsymbol{\theta}$:

$$s_\theta(i) = \begin{cases} \text{sign}(\theta(i)), & i \in \Gamma \\ 0, & i \notin \Gamma \end{cases}.$$

We consider two cases: $\mathbf{s}_\theta^\top \mathbf{z} \neq 0$ and $\mathbf{s}_\theta^\top \mathbf{z} = 0$.

First, suppose that $\mathbf{s}_\theta^\top \mathbf{z} \neq 0$. Without loss of generality, we can assume that $\mathbf{s}_\theta^\top \mathbf{z} < 0$, as otherwise we can just replace \mathbf{z} with $-\mathbf{z}$ (since both of these are supported on Γ and are in the nullspace of \mathbf{A}). Notice that we can write

$$\|\boldsymbol{\theta}\|_1 = \sum_{i=1}^{d+1} |\theta(i)| = \sum_{i=1}^{d+1} \text{sign}(\theta(i))\theta(i)$$

For $\epsilon > 0$ small enough, it is a fact that

$$\text{sign}(\theta(i) + \epsilon z(i)) = \text{sign}(\theta(i)), \quad \text{for all } i \in \Gamma,$$

and so

$$\begin{aligned}
 \|\boldsymbol{\theta} + \epsilon \mathbf{z}\|_1 &= \sum_{i=1}^{d+1} \text{sign}(\theta(i) + \epsilon z(i))(\theta(i) + \epsilon z(i)) \\
 &= \sum_{i=1}^{d+1} \text{sign}(\theta(i))(\theta(i) + \epsilon z(i)) \\
 &= \|\boldsymbol{\theta}\|_1 + \epsilon \mathbf{s}_\theta^T \mathbf{z} \\
 &< \|\boldsymbol{\theta}\|_1.
 \end{aligned}$$

Thus, taking $\boldsymbol{\theta}' = \boldsymbol{\theta} + \epsilon \mathbf{z}$ for a small enough value of $\epsilon > 0$ gives a vector with a smaller functional value, so $\boldsymbol{\theta}$ cannot be a solution to (LASSO).

Now suppose that $\mathbf{s}_\theta^T \mathbf{z} = 0$. Then

$$\|\boldsymbol{\theta}\|_1 = \|\boldsymbol{\theta} + \epsilon \mathbf{z}\|_1$$

as long as $\boldsymbol{\theta} + \epsilon \mathbf{z}$ remains non-zero on Γ . But there must be at least one $i \in \Gamma$ such that

$$\text{sign}(\theta(i)) \neq \text{sign}(z(i)),$$

otherwise the inner product could not be equal to zero. Let i' be the index that obeys the condition above such that $\theta(i)$ is smallest relative to $z(i)$:

$$i' = \arg \min_{i \in \Gamma} \left\{ \frac{|\theta(i)|}{|z(i)|} : \text{sign}(\theta(i)) \neq \text{sign}(z(i)) \right\}.$$

Then

$$\boldsymbol{\theta}' = \boldsymbol{\theta} + \frac{|\theta(i')|}{|z(i')|} \mathbf{z},$$

will be exactly equal to zero at i' and will still be non-zero outside of Γ . Thus

$$\text{nnz}(\boldsymbol{\theta}') < \text{nnz}(\boldsymbol{\theta}),$$

while

$$\frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}'\|_2^2 + \lambda\|\boldsymbol{\theta}'\|_1 = \frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda\|\boldsymbol{\theta}\|_1.$$

The LASSO as a quadratic program

Unlike standard least-squares, or least-squares with Tikhonov regularization, the solution to the LASSO does not have a closed form. We can, however, write it as a convex quadratic program with linear inequality constraints. This puts it in the same class of optimization program as SVMs.

The main idea is to introduce slack variables that allow us to re-write the ℓ_1 norm, which is piecewise linear, as a linear function subject to linear constraints. In particular, the solution to the LASSO is exactly the same as the solution to

$$\begin{aligned} \underset{\boldsymbol{\theta}, \mathbf{u}}{\text{minimize}} \quad & \frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda \sum_{i=1}^{d+1} u(i) \\ \text{subject to} \quad & -u(i) \leq \theta(i) \leq u(i), \quad i = 1, \dots, d+1. \end{aligned}$$

This is the same as solving

$$\begin{aligned} \underset{\boldsymbol{\theta}, \mathbf{u}}{\text{minimize}} \quad & \frac{1}{2}\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda \sum_{i=1}^{d+1} u(i) \\ \text{subject to} \quad & \boldsymbol{\theta} - \mathbf{u} \leq \mathbf{0} \\ & -\boldsymbol{\theta} - \mathbf{u} \leq \mathbf{0}, \end{aligned}$$

which is the same as solving

$$\underset{z \in \mathbb{R}^{2d+2}}{\text{minimize}} \quad \frac{1}{2} z^T \mathbf{P} z + \mathbf{c}^T z \quad \text{subject to} \quad \mathbf{R} z \leq \mathbf{0},$$

where

$$z = \begin{bmatrix} \boldsymbol{\theta} \\ \mathbf{u} \end{bmatrix}, \quad \mathbf{P} = \begin{bmatrix} \mathbf{A}^T \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} -\mathbf{A}^T \mathbf{y} \\ \lambda \mathbf{1} \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} \mathbf{I} & -\mathbf{I} \\ -\mathbf{I} & -\mathbf{I} \end{bmatrix}.$$