

Least-Squares in (∞ -dimensional) Hilbert Space

We have seen that introducing a (finite dimensional) basis allows us to turn the regression problem into a finite-dimensional linear algebra problem. In this section, we will see that even without a basis, we can solve least-squares problems in ∞ -dimensional Hilbert space by recasting them as finite-dimensional linear algebra problems.

To start, let's look back at our solution to the ridge regression problem when \mathbf{A} is underdetermined (more columns than rows, $M < N$). We solve

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta \|\mathbf{x}\|_2^2,$$

which we can also rewrite as

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \sum_{m=1}^M |y_m - \mathbf{a}_m^T \mathbf{x}|^2 + \delta \|\mathbf{x}\|_2^2,$$

where the $\mathbf{a}_1, \dots, \mathbf{a}_M \in \mathbb{R}^N$ are the rows of \mathbf{A} (after we transpose them):

$$\mathbf{A} = \begin{bmatrix} -\mathbf{a}_1^T - \\ -\mathbf{a}_2^T - \\ \vdots \\ -\mathbf{a}_M^T - \end{bmatrix}.$$

We have seen that in this case the solution is

$$\hat{\mathbf{x}} = \left(\mathbf{A}^T \mathbf{A} + \delta \mathbf{I} \right)^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}^T \left(\mathbf{A} \mathbf{A}^T + \delta \mathbf{I} \right)^{-1} \mathbf{y}$$

We explicitly point out two facts about the solution:

1. $\hat{\mathbf{x}}$ is in the *row space* (the linear span of the rows) of \mathbf{A}

$$\hat{\mathbf{x}} = \mathbf{A}^T \hat{\boldsymbol{\alpha}} = \sum_{m=1}^M \hat{\alpha}_m \mathbf{a}_m.$$

2. The coefficients $\hat{\boldsymbol{\alpha}}$ are computed by solving the symmetric positive definite system of equations

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{y}, \quad \mathbf{K} = \mathbf{A} \mathbf{A}^T.$$

The $M \times M$ matrix \mathbf{K} is formed by taking all the inner products between the different rows of \mathbf{A} :

$$\mathbf{K} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{a}_1 & \mathbf{a}_1^T \mathbf{a}_2 & \cdots & \mathbf{a}_1^T \mathbf{a}_M \\ \mathbf{a}_2^T \mathbf{a}_1 & \mathbf{a}_2^T \mathbf{a}_2 & \cdots & \mathbf{a}_2^T \mathbf{a}_M \\ \vdots & & \ddots & \vdots \\ \mathbf{a}_M^T \mathbf{a}_1 & \cdots & & \mathbf{a}_M^T \mathbf{a}_M \end{bmatrix}.$$

We will see below that these two facts extend to the analogous problem in an ∞ -dimensional Hilbert space ... this is convenient in that it allows us to solve ∞ -dimensional optimization problems with a finite amount of computational effort.

Let \mathcal{S} be a Hilbert space with inner product $\langle \cdot, \cdot \rangle_S$. For the purposes of discussion, we will think of the elements of \mathcal{S} as containing functions that map $\mathbb{R}^D \rightarrow \mathbb{R}$, but the math we use below applies to any abstract Hilbert space. Let \mathbf{f} be a function in \mathcal{S} that we are trying to estimate from observations of a series of inner products against fixed $\mathbf{a}_1, \dots, \mathbf{a}_M \in \mathcal{S}$:

$$\begin{aligned} y_1 &= \langle \mathbf{f}, \mathbf{a}_1 \rangle_S + \text{noise} \\ y_2 &= \langle \mathbf{f}, \mathbf{a}_2 \rangle_S + \text{noise} \\ &\vdots \\ y_M &= \langle \mathbf{f}, \mathbf{a}_M \rangle_S + \text{noise} \end{aligned}$$

This is analagous to observing $\mathbf{y} = \mathbf{A} \mathbf{x} + \text{noise}$ in the finite dimensional case (where $\mathbf{x} \in \mathbb{R}^N$), where each entry in \mathbf{y} is being modeled

as an inner product between a row in \mathbf{A} and \mathbf{x} . Here, it's sort of like we are observing $\mathbf{y} \in \mathbb{R}^M$ through a “matrix” that has M rows, but each row is a function in \mathcal{S} instead of a vector in \mathbb{R}^N .

Suppose now that we estimate \mathbf{f} by solving the following least-squares problem in \mathcal{S} :

$$\underset{\mathbf{f} \in \mathcal{S}}{\text{minimize}} \quad \sum_{m=1}^M |y_m - \langle \mathbf{f}, \mathbf{a}_m \rangle_S|^2 + \delta \|\mathbf{f}\|_S^2. \quad (1)$$

This is an optimization program in a (possibly) infinite dimensional Hilbert space. Even if we discretize the search space using an orthonobasis (or any basis), there will be an infinite number of expansion coefficients to solve for. However, the following result (which is almost an immediate consequence of our work on linear approximation a little while ago) shows us that we can at least specify the solution by solving an $M \times M$ system of equations.

Representer Theorem (least-squares version):

The solution to (1) is given by

$$\hat{\mathbf{f}} = \sum_{m=1}^M \hat{\alpha}_m \mathbf{a}_m,$$

where

$$\hat{\alpha} = (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{y}, \quad \mathbf{K} = \begin{bmatrix} \langle \mathbf{a}_1, \mathbf{a}_1 \rangle_S & \langle \mathbf{a}_2, \mathbf{a}_1 \rangle_S & \cdots & \langle \mathbf{a}_M, \mathbf{a}_1 \rangle_S \\ \langle \mathbf{a}_1, \mathbf{a}_2 \rangle_S & \langle \mathbf{a}_2, \mathbf{a}_2 \rangle_S & \cdots & \langle \mathbf{a}_M, \mathbf{a}_2 \rangle_S \\ \vdots & & \ddots & \vdots \\ \langle \mathbf{a}_1, \mathbf{a}_M \rangle_S & \cdots & & \langle \mathbf{a}_M, \mathbf{a}_M \rangle_S \end{bmatrix}.$$

Proof. We use the notation

$$L(\mathbf{f}) = \sum_{m=1}^M |y_m - \langle \mathbf{f}, \mathbf{a}_m \rangle_S|^2,$$

and so we are trying to solve

$$\underset{\mathbf{f}}{\text{minimize}} L(\mathbf{f}) + \delta \|\mathbf{f}\|_S^2.$$

Let $\mathcal{A} = \text{Span}(\{\mathbf{a}_1, \dots, \mathbf{a}_M\})$ be the subspace spanned by the \mathbf{a}_m . For any candidate function $\mathbf{g} \in \mathcal{S}$, we can write

$$\mathbf{g} = \mathbf{g}_A + \mathbf{g}_\perp,$$

where \mathbf{g}_A is the closest point in \mathcal{A} to \mathbf{g} and $\mathbf{g}_\perp = \mathbf{g} - \mathbf{g}_A$ is orthogonal to every vector in \mathcal{A} ; in particular

$$\langle \mathbf{g}_\perp, \mathbf{a}_m \rangle_S = 0, \quad m = 1, \dots, M.$$

Then

$$\begin{aligned} L(\mathbf{g}) &= \sum_{m=1}^M |y_m - \langle \mathbf{g}_A + \mathbf{g}_\perp, \mathbf{a}_m \rangle_S|^2 \\ &= \sum_{m=1}^M |y_m - \langle \mathbf{g}_A, \mathbf{a}_m \rangle_S|^2 \quad (\text{since } \langle \mathbf{g}_\perp, \mathbf{a}_m \rangle_S = 0) \\ &= L(\mathbf{g}_A), \end{aligned}$$

and

$$\begin{aligned} \|\mathbf{g}\|_S^2 &= \|\mathbf{g}_A\|_S^2 + \|\mathbf{g}_\perp\|_S^2 \quad (\text{Pythagorean thm.}) \\ &\geq \|\mathbf{g}_A\|_S^2. \end{aligned}$$

Thus

$$L(\mathbf{g}) + \delta \|\mathbf{g}\|_S^2 \geq L(\mathbf{g}_A) + \delta \|\mathbf{g}_A\|_S^2.$$

So for every $\mathbf{g} \in \mathcal{S}$, there is a corresponding member of \mathcal{A} that makes the functional we are trying to minimize at least as small. Thus at least one solution to (1) must be in \mathcal{A} , and we can write

$$\hat{\mathbf{f}} = \sum_{m=1}^M \hat{\alpha}_m \mathbf{a}_m, \quad (2)$$

for some $\hat{\alpha}_1, \dots, \hat{\alpha}_M$. Note that if the \mathbf{a}_m are not all linearly independent, then there might be multiple $\hat{\alpha}$ that give the same $\hat{\mathbf{f}}$, but this is not a problem; we only need to find one of them. No matter what, we still have that $\hat{\mathbf{f}}$ is the **unique** solution to (1), there just might be multiple ways to write that solution.

We pause here to specifically point out something that will make everything a little easier below. If \mathbf{v} and \mathbf{u} are elements in \mathcal{A} with

$$\mathbf{u} = \sum_{m=1}^M c_m \mathbf{a}_m, \quad \text{and} \quad \mathbf{v} = \sum_{m=1}^M d_m \mathbf{a}_m,$$

then

$$\begin{aligned} \langle \mathbf{u}, \mathbf{v} \rangle_S &= \left\langle \sum_{m=1}^M c_m \mathbf{a}_m, \sum_{\ell=1}^M d_\ell \mathbf{a}_\ell \right\rangle_S \\ &= \sum_{m=1}^M \sum_{\ell=1}^M c_m d_\ell \langle \mathbf{a}_m, \mathbf{a}_\ell \rangle_S \\ &= \mathbf{d}^T \mathbf{K} \mathbf{c}, \end{aligned}$$

where \mathbf{K} is the same matrix given in the statement of the theorem — it is similar to the Gram matrix for the $\{\mathbf{a}_m\}$, though there is no

guarantee here that the $\{\mathbf{a}_m\}$ are linearly independent, and so \mathbf{K} might not be invertible.

To find the best $\boldsymbol{\alpha}$ in (2), we solve

$$\begin{aligned}
& \underset{\boldsymbol{\alpha} \in \mathbb{R}^M}{\text{minimize}} \sum_{m=1}^M \left| y_m - \sum_{\ell=1}^M \alpha_\ell \langle \mathbf{a}_\ell, \mathbf{a}_m \rangle \right|^2 + \delta \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \\
& \quad \downarrow \\
& \underset{\boldsymbol{\alpha} \in \mathbb{R}^M}{\text{minimize}} \sum_{m=1}^M |y_m - (\mathbf{K} \boldsymbol{\alpha})_m|^2 + \delta \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \\
& \quad \downarrow \\
& \underset{\boldsymbol{\alpha} \in \mathbb{R}^M}{\text{minimize}} \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \delta \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}.
\end{aligned}$$

Taking the gradient of the functional above and setting it equal to zero¹ tells us that any solution $\hat{\boldsymbol{\alpha}}$ must obey

$$\mathbf{K}^T(\mathbf{K} \hat{\boldsymbol{\alpha}} - \mathbf{y}) + \delta \mathbf{K} \hat{\boldsymbol{\alpha}} = \mathbf{0}. \quad (3)$$

The matrix \mathbf{K} is symmetric ($\mathbf{K}^T = \mathbf{K}$), and so any $\hat{\boldsymbol{\alpha}}$ that obeys

$$(\mathbf{K} + \delta \mathbf{I}) \hat{\boldsymbol{\alpha}} = \mathbf{y},$$

will also obey (3). Since the matrix $\mathbf{K} + \delta \mathbf{I}$ is always invertible, we can take $\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{y}$ and the theorem is established. ■

¹The gradient (first derivative) being zero is always a necessary condition for minimizer of a functional with multiple arguments. In this case, it happens to be sufficient as well since the Hessian matrix (second derivative) is positive semi-definite, meaning the functional is convex.