

## II. Regression using Least Squares

The fundamental problem in supervised machine learning is to fit a function to a series of point evaluations. We are given data

$$(\mathbf{t}_1, y_1), \dots, (\mathbf{t}_M, y_m), \quad \text{with } \mathbf{t}_m \in \mathbb{R}^D \text{ and } y_m \in \mathbb{R},$$

and want to find a function  $f : \mathbb{R}^D \rightarrow \mathbb{R}$  such that

$$f(\mathbf{t}_m) \approx y_m, \quad m = 1, \dots, M. \tag{1}$$

When the  $y_m$  (and the range of  $f$ ) are continuous-valued, the process of fitting such an  $f$  is called **regression**. This is a terrible name whose English meaning is only tangentially related to this problem. But the terminology has become ingrained through 130+ years of use, so we are stuck with it.

We need two key ingredients to putting the task in (1) on firm mathematical ground, a model function class and a loss function.

**Function class.** For the problem (1) to make sense, our search needs to be limited to a class of functions  $\mathcal{F}$ . For example, we might restrict  $f$  to be polynomial, or twice differentiable, etc. Choosing this  $\mathcal{F}$  is a modeling problem, as it basically encodes what we believe to be true (or reasonably close to true) about the function we are trying to discover. It also has to be something we can compute with. Examples below will take  $\mathcal{F}$  to be an appropriately chosen Hilbert space, or a subspace of a Hilbert space; these will both lead to algorithms that lean heavily on linear algebra.

In some cases, this  $\mathcal{F}$  is chosen implicitly. This is the case in modern multi-layer neural networks, where the structure for computing the function is defined, but it is hard to describe exactly the class of functions they can represent.

With the choice of  $\mathcal{F}$  fixed, the problem (1) now becomes

$$\text{find } f \in \mathcal{F} \text{ such that } f(\mathbf{t}_m) \approx y_m, \quad m = 1, \dots, M. \tag{2}$$

**Loss function.** This penalizes the deviation of the  $f(\mathbf{t}_m)$  from the  $y_m$ . Given a loss function  $\ell(\cdot, \cdot) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$  that quantifies the penalty for the deviation at a single sample location, we can assign a (positive) numeric score to the performance of every candidate  $f$  by summing over all the samples. This allows us to write (2) more precisely as an optimization problem:

$$\underset{f \in \mathcal{F}}{\text{minimize}} \sum_{m=1}^M \ell(f(\mathbf{t}_m), y_m) \quad (3)$$

There are again many loss functions you might consider, and depending on the context, some might be more natural than others. We will focus almost all of our efforts on the **squared loss**  $\ell(u, v) = |u - v|^2$ . Then (3) becomes

$$\underset{f \in \mathcal{F}}{\text{minimize}} \sum_{m=1}^M |y_m - f(\mathbf{t}_m)|^2. \quad (4)$$

We will see below that this choice, coupled with a subspace or Hilbert space model for  $\mathcal{F}$ , allows us to completely solve the regression problem using linear algebra.

## Linear regression

Perhaps the most classical choice for the function class  $\mathcal{F}$  is that it contains all linear functions on  $\mathbb{R}^D$ . A function  $f : \mathbb{R}^D \rightarrow \mathbb{R}$  is **linear** if

$$f(\alpha \mathbf{t}_1 + \beta \mathbf{t}_2) = \alpha f(\mathbf{t}_1) + \beta f(\mathbf{t}_2), \quad \text{for all } \alpha, \beta \in \mathbb{R}, \mathbf{t}_1, \mathbf{t}_2 \in \mathbb{R}^D.$$

It is a fact<sup>1</sup> that every linear functional on  $\mathbb{R}^D$  is uniquely represented by a vector  $\mathbf{c}_f \in \mathbb{R}^D$ , where

$$f(\mathbf{t}) = \langle \mathbf{t}, \mathbf{w}_f \rangle = \mathbf{t}^\top \mathbf{w}_f.$$

So given the  $\{(\mathbf{t}_m, y_m)\}$ , we want to find  $\mathbf{w}$  such that

$$y_m \approx \mathbf{t}_m^\top \mathbf{w},$$

and using the squared-loss to measure the mismatch, we have the optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^D}{\text{minimize}} \quad \sum_{m=1}^M |y_m - \mathbf{t}_m^\top \mathbf{w}|^2.$$

If we stacking up the  $\mathbf{t}_m$  as rows in a  $M \times D$  matrix  $\mathbf{A}$ , the sum in the optimization program above becomes  $\|\mathbf{y} - \mathbf{A}\mathbf{w}\|_2^2$ , where

$$\mathbf{A} = \begin{bmatrix} - & \mathbf{t}_1^\top & - \\ - & \mathbf{t}_2^\top & - \\ & \vdots & \\ - & \mathbf{t}_M^\top & - \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix}.$$

We can then find the best linear function on  $\mathbb{R}^D$  by solving the finite-dimensional least-squares problem

$$\underset{\mathbf{w} \in \mathbb{R}^D}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{w}\|_2^2.$$

Given a solution  $\hat{\mathbf{w}}$  to the above, we have the best linear function  $\hat{f}$ , where

$$\hat{f}(\mathbf{t}) = \mathbf{t}^\top \hat{\mathbf{w}} = \hat{w}_1 t_1 + \hat{w}_2 t_2 + \cdots + \hat{w}_D t_D.$$

---

<sup>1</sup>The name by which this fact goes is the *Riesz Representation Theorem*, which we will encounter a little later in the course.

It is often the case that we want to also want to add a constant offset to the above, and find an **affine** function of the form

$$f(\mathbf{t}) = w_0 + w_1 t_1 + \cdots + w_D t_D.$$

This is done simply by adding a column of all ones to  $\mathbf{A}$ :

$$\mathbf{A}' = \begin{bmatrix} 1 & - & \mathbf{t}_1^T & - \\ 1 & - & \mathbf{t}_2^T & - \\ \vdots & & \vdots & \\ 1 & - & \mathbf{t}_M^T & - \end{bmatrix},$$

and then solving

$$\underset{\mathbf{w}' \in \mathbb{R}^{D+1}}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}'\mathbf{w}'\|_2^2. \quad (5)$$

## Nonlinear regression using a basis

It is easy to generalize what is in the previous section to fit nonlinear functions to data. What is interesting is that we still end up with a very similar linear least-squares problem.

We do this by letting  $\mathcal{F}$  be a finite-dimensional subspace (of a Hilbert space) spanned by a set of specified basis functions. Given a set of building blocks  $\boldsymbol{\psi}_n : \mathbb{R}^D \rightarrow \mathbb{R}$  (i.e. basis functions), the model implicit here is that our target function  $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}$  is (at least approximately) in the subspace spanned by  $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_N$ , that is, there exists  $\{x_n\}$  such that

$$f(\mathbf{t}) = \sum_{n=1}^N x_n \psi_n(\mathbf{t}).$$

Fitting a function in this space  $\mathcal{F}$  is the same as fitting a  $\mathbf{x} \in \mathbb{R}^N$  such that

$$y_1 \approx \sum_{n=1}^N x_n \psi_n(\mathbf{t}_1), \quad y_2 \approx \sum_{n=1}^N x_n \psi_n(\mathbf{t}_2), \quad \dots \quad y_M \approx \sum_{n=1}^N x_n \psi_n(\mathbf{t}_M).$$

With the least-squares loss, we want to solve

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \sum_{m=1}^M \left| y_m - \sum_{n=1}^N x_n \psi_n(\mathbf{t}_m) \right|^2 = \sum_{m=1}^M \left| y_m - \mathbf{\Psi}(\mathbf{t}_m)^T \mathbf{x} \right|^2, \quad (6)$$

where  $\mathbf{\Psi}(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^N$  is

$$\mathbf{\Psi}(\mathbf{t}) = \begin{bmatrix} \psi_1(\mathbf{t}) \\ \psi_2(\mathbf{t}) \\ \vdots \\ \psi_N(\mathbf{t}) \end{bmatrix}.$$

Constructing the  $M \times N$  matrix  $\mathbf{A}$  as

$$\mathbf{A} = \begin{bmatrix} - & \mathbf{\Psi}(\mathbf{t}_1)^T & - \\ - & \mathbf{\Psi}(\mathbf{t}_2)^T & - \\ & \vdots & \\ - & \mathbf{\Psi}(\mathbf{t}_M)^T & - \end{bmatrix} = \begin{bmatrix} \psi_1(\mathbf{t}_1) & \psi_2(\mathbf{t}_1) & \cdots & \psi_N(\mathbf{t}_1) \\ \psi_1(\mathbf{t}_2) & \psi_2(\mathbf{t}_2) & \cdots & \psi_N(\mathbf{t}_2) \\ \vdots & & \ddots & \\ \psi_1(\mathbf{t}_M) & \psi_2(\mathbf{t}_M) & \cdots & \psi_N(\mathbf{t}_M) \end{bmatrix}, \quad (7)$$

the optimization program (4) becomes

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2. \quad (8)$$

From a solution  $\hat{\mathbf{x}}$  to the program above, we can synthesize the solution to (6) as

$$\hat{f}(\mathbf{t}) = \sum_{n=1}^N \hat{x}_n \psi_n(\mathbf{t}).$$

Even though we are ultimately recovering a nonlinear function of a continuous variable, introducing a basis allows us to put the nonlinear regression problem into the exact same computational framework as linear regression.

## Example

Play around with the code in `regression_examples.m`. Here, we take  $M$  (noisy) samples of an underlying function, and then can perform regression using  $N$  basis functions .... code for polynomials, “bumps”, and Fourier (trigonometric polynomial) are all included.

Notice that increasing  $N$  makes the class richer, and corresponds to adding columns in  $\mathbf{A}$ . This increases our ability to match the samples, but also increases the risk of “over fitting” the model.

# The least-squares problem

We have seen that in both the linear regression case (5) and the nonlinear regression case (5), the problem reduces to same form. We will look at problems like this multiple times in this course, and it is worth starting to study this kind of problem from the perspective of linear algebra.

We start with the following fundamental result:

Any solution  $\hat{\mathbf{x}}$  to

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \quad (9)$$

must obey the **normal equations**

$$\mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{y}. \quad (10)$$

It is a fact<sup>2</sup> that  $g(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$  is a convex, differentiable function on all of  $\mathbb{R}^N$ . Thus a necessary condition for  $\hat{\mathbf{x}}$  to be minimizer of  $g(\mathbf{x})$  is that the gradient (with respect to  $\mathbf{x}$ ) is equal to zero at  $\hat{\mathbf{x}}$ ,  $\nabla g(\hat{\mathbf{x}}) = \mathbf{0}$ . We have

$$\begin{aligned} \nabla (\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2) &= \nabla (\|\mathbf{y}\|_2^2 - 2\mathbf{y}^T \mathbf{A}\mathbf{x} + \|\mathbf{A}\mathbf{x}\|_2^2) \\ &= -2\mathbf{A}^T \mathbf{y} + 2\mathbf{A}^T \mathbf{A}\mathbf{x}, \end{aligned}$$

which is  $\mathbf{0}$  when (10) holds.

There are several natural questions at this point. First, is there any guarantee that (10) even has a solution? Are there conditions under

---

<sup>2</sup>Prove this at home!



which the solution is guaranteed to be unique or non-unique? And if the solution is non-unique, what should we do?

We will answer these questions using the following two facts from linear algebra. For any  $M \times N$  matrix  $\mathbf{A}$ ,

- $\text{Null}(\mathbf{A}^T \mathbf{A}) = \text{Null}(\mathbf{A})$ , and
- $\text{Col}(\mathbf{A}^T \mathbf{A}) = \text{Row}(\mathbf{A})$ .

Proof of these facts can be found in the Technical Details section.

It is now easy to argue that:

1. The system (10) always has a solution, no matter what  $\mathbf{A}$  and  $\mathbf{y}$  are. This follows immediately from the fact that

$$\mathbf{A}^T \mathbf{y} \in \text{Row}(\mathbf{A}) = \text{Col}(\mathbf{A}^T \mathbf{A}).$$

This means that there is always at least one minimizer of (9).

2. If  $\text{rank}(\mathbf{A}) = N$ , then (10) has the unique solution, and so (9) has a unique minimizer. This unique minimizer is simply

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}.$$

3. If  $\text{rank}(\mathbf{A}) < N$ , then there are an infinite number of solutions to (9). In this case,  $\mathbf{A}$  has a non-trivial null space, so if  $\hat{\mathbf{x}}$  is a solution to (9), so is  $\hat{\mathbf{x}} + \mathbf{v}$  for all  $\mathbf{v} \in \text{Null}(\mathbf{A})$ , as

$$\|\mathbf{y} - \mathbf{A}(\hat{\mathbf{x}} + \mathbf{v})\|_2^2 = \|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}} - \mathbf{A}\mathbf{v}\|_2^2 = \|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_2^2,$$

since  $\mathbf{A}\mathbf{v} = \mathbf{0}$ .

4. If  $\text{rank}(\mathbf{A}) = M$ , then there exists at least one  $\hat{\mathbf{x}}$  such that  $\|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_2^2 = 0$ , that is  $\mathbf{A}\hat{\mathbf{x}} = \mathbf{y}$ . If in addition  $M < N$ , then there will be an infinity of such solutions.

## Minimum $\ell_2$ solutions

When there are an infinite number of solutions, we need a way to pick one of them. For this, modeling will again come into play ... we need some kind of value system other than the objective in (9) to judge which of the solutions is best. There are many, many ways this can be done; we will look two variations of one idea below.

One principle is to choose the solution that is the smallest (i.e. closest to the origin). In many applications this is justified using some kind of *minimum energy* principle — the sum of squares being a proxy for the amount of resources it takes to implement a solution. It can also be thought of as a kind of Occam's razor, in that the solution should not be any bigger than it needs to be.

The minimum energy least-squares solution is now the solution to the optimization program

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}. \quad (11)$$

It happens that there is now always a unique solution to this program, we will see why this is true when we revisit this problem in a later section of the course (and look at it through the lens of the singular value decomposition. For now, we will look at the case where  $\mathbf{A}$  has full row rank,  $\text{rank}(\mathbf{A}) = M$ . In this case,  $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}$  if and only if  $\mathbf{A} \mathbf{x} = \mathbf{y}$ , so (11) simplifies to

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad \mathbf{A} \mathbf{x} = \mathbf{y}. \quad (12)$$

The first thing to realize (and this holds in the general case (11) as well), is that the solution to the above will be in  $\text{Row}(\mathbf{A})$ . We

know that the row and null spaces are orthogonal complements of one another, and so every  $\mathbf{x} \in \mathbb{R}^N$  can be written as

$$\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2, \quad \text{where } \mathbf{x}_1 \in \text{Row}(\mathbf{A}), \mathbf{x}_2 \in \text{Null}(\mathbf{A}),$$

and of course  $\mathbf{x}_1^T \mathbf{x}_2 = 0$ . We can recast the optimization in (12) as a search over  $\mathbf{x}_1$  and  $\mathbf{x}_2$ :

$$\begin{array}{ll} \underset{\substack{\mathbf{x}_1 \in \text{Row}(\mathbf{A}) \\ \mathbf{x}_2 \in \text{Null}(\mathbf{A})}}{\text{minimize}} & \|\mathbf{x}_1 + \mathbf{x}_2\|_2^2 \quad \text{subject to} \quad \mathbf{A}(\mathbf{x}_1 + \mathbf{x}_2) = \mathbf{y}. \end{array}$$

Since  $\mathbf{A}\mathbf{x}_2 = \mathbf{0}$  and  $\|\mathbf{x}_1 + \mathbf{x}_2\|_2^2 = \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2$  (Pythagorous), the above simplifies to

$$\begin{array}{ll} \underset{\substack{\mathbf{x}_1 \in \text{Row}(\mathbf{A}) \\ \mathbf{x}_2 \in \text{Null}(\mathbf{A})}}{\text{minimize}} & \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 \quad \text{subject to} \quad \mathbf{A}\mathbf{x}_1 = \mathbf{y}. \end{array}$$

Since now  $\mathbf{x}_2$  does not appear in the constraints, we see that the minimizer will have  $\mathbf{x}_2 = \mathbf{0}$ , i.e. the solution lies entirely in  $\text{Row}(\mathbf{A})$ .

When  $\text{rank}(\mathbf{A}) = M$ , we can use the fact that  $\hat{\mathbf{x}} \in \text{Row}(\mathbf{A})$  to derive a closed-form solution. We know that there exists a  $\hat{\mathbf{v}}$  such that

$$\hat{\mathbf{x}} = \mathbf{A}^T \hat{\mathbf{v}}, \quad \text{and} \quad \mathbf{A}\mathbf{A}^T \hat{\mathbf{v}} = \mathbf{y}.$$

Since  $\text{rank}(\mathbf{A}) = M$ , we know that the  $M \times M$  matrix  $\mathbf{A}\mathbf{A}^T$  is invertible, and so there is exactly one  $\hat{\mathbf{v}}$  that obeys the second condition above, namely  $\hat{\mathbf{v}} = (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{y}$ . Thus the solution to (12) is

$$\hat{\mathbf{x}} = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{y}.$$

We will derive the solution to this problem for  $\text{rank}(\mathbf{A}) < M$  in a few weeks when we talk about the singular value decomposition.

## Regularization

The problem with solving (11) when  $\mathbf{A}$  has a non-trivial null space is not just that there are an infinity of solutions, it is also that the space of solutions is unbounded — you can add a vector from the null space of arbitrary size and not change the functional. Even when  $\mathbf{A}$  technically has full column rank, if it is poorly conditioned, meaning that

$$\|(\mathbf{A}^T \mathbf{A})^{-1}(\mathbf{x} - \mathbf{z})\|_2 \gg \|\mathbf{x} - \mathbf{z}\|_2, \quad \text{for some } \mathbf{x}, \mathbf{z} \in \mathbb{R}^N,$$

then a small change in  $\mathbf{y}$  can amount to a massive change in  $\hat{\mathbf{x}}$ .

We will look much more carefully at describing how well-conditioned the least-squares problem is when we talk about the SVD. For now, know that there is a way to temper this behavior by penalizing the size of the solution. In place of (11), we solve

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta \|\mathbf{x}\|_2^2, \quad (13)$$

for some  $\delta \geq 0$ . It is similar in spirit to the minimum norm problem in (12) in that we are favoring solutions close to the origin, but it can be applied for any  $\mathbf{A}$ , no matter the rank. It also gives us a little more flexibility in the model, as we can treat  $\delta$  as a “knob” that sets the trade-off between how closely we want  $\mathbf{A}\mathbf{x}$  to match  $\mathbf{y}$ , and how large  $\mathbf{x}$  can be. Note that when  $\text{rank}(\mathbf{A}) = M$ , the solution to (13) goes to the solution to (12) as  $\delta \rightarrow 0$ .

The analog of the normal equations for (13) is that a solution  $\hat{\mathbf{x}}$  must obey

$$(\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{y}. \quad (14)$$

However, the matrix  $\mathbf{A}^T \mathbf{A} + \delta \mathbf{I}$  is always invertible for  $\delta > 0$  no matter what  $\mathbf{A}$  is (we will see proof of this later in the course). Thus

the unique solution to (13) is

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

It also turns out that we can write this as

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

You can verify this by plugging the expression above into the left hand side of (14).

From a high level perspective, solving (13) in place of (9) recognizes that making  $\mathbf{Ax}$  as close to  $\mathbf{y}$  as possible shouldn't be the only thing we are interested in. Models (which ultimately is where  $\mathbf{A}$  comes from) are rarely perfect, and so pulling out all the stops to push  $\mathbf{Ax}$  just a little closer to  $\mathbf{y}$  can be self-defeating.

In the context of regression, using (13) in place of (11) is often referred to as **ridge regression**.

## Ridge regression using a basis

Let's return to our fundamental problem of fitting a function to data (regression). Above, we saw that once a basis was introduced, we could set this up as a least-squares problem (equations (7) and (8) above), and then we saw one way to make this least-squares problem better posed using regularization.

Another variation is to impose the regularization in the function space. Suppose that we are looking for an  $\mathbf{f}$  in a Hilbert space  $\mathcal{S}$  that match observed data  $\{(\mathbf{t}_m, y_m)\}$ . The associated regularized

least-squares problem is

$$\underset{f \in \mathcal{S}}{\text{minimize}} \sum_{m=1}^M |y_m - f(\mathbf{t}_m)|^2 + \delta \|\mathbf{f}\|_S^2. \quad (15)$$

Note that we are using the Hilbert space norm to penalize the size of the function. If we again introduce a basis, modeling the target function as being in the span of  $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_N$ , then we can rewrite this as

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta \left\| \sum_{n=1}^N x_n \boldsymbol{\psi}_n \right\|_S^2,$$

where  $\mathbf{A}$  is constructed as in (7). Expanding the term on the right above gives us

$$\begin{aligned} \left\| \sum_{n=1}^N x_n \boldsymbol{\psi}_n \right\|_S^2 &= \left\langle \sum_{n=1}^N x_n \boldsymbol{\psi}_n, \sum_{k=1}^N x_k \boldsymbol{\psi}_k \right\rangle_S \\ &= \sum_{n=1}^N \sum_{k=1}^N x_n x_k \langle \boldsymbol{\psi}_n, \boldsymbol{\psi}_k \rangle_S \\ &= \mathbf{x}^T \mathbf{G} \mathbf{x}, \end{aligned}$$

where  $\mathbf{G}$  is the Gram matrix for the basis,  $G_{n,k} = \langle \boldsymbol{\psi}_n, \boldsymbol{\psi}_k \rangle_S$ . Thus the optimization (15) in Hilbert space becomes the finite-dimensional problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta \mathbf{x}^T \mathbf{G} \mathbf{x}, \quad (16)$$

which has closed form solution<sup>3</sup>

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A} + \delta \mathbf{G})^{-1} \mathbf{A}^T \mathbf{y}.$$

---

<sup>3</sup>Actually, there is no guarantee here that  $\mathbf{A}^T \mathbf{A} + \delta \mathbf{G}$  is invertible for all  $\delta > 0$ , but we will assume that it is.

We then synthesize the solution as before,

$$\hat{f}(\mathbf{t}) = \sum_{n=1}^N x_n \psi_n(\mathbf{t}).$$

Of course, when the  $\{\boldsymbol{\psi}_n\}$  are orthonormal, then  $\mathbf{G} = \mathbf{I}$ , and (16) is equivalent to the standard regularized least-square problem (13). When the basis is not orthonormal, then solving (16) and (13) will in general produce different weights (and hence different synthesized functions). But in practice, if the basis is reasonable (meaning that  $\mathbf{G}$  is well-conditioned), then the difference will be minor.

## Technical Details: Two easy linear algebra facts

**Proposition.** Let  $\mathbf{A}$  be an  $M \times N$  matrix. Then  $\text{Null}(\mathbf{A}^T \mathbf{A}) = \text{Null}(\mathbf{A})$ .

**Proof.** It is clear that  $\mathbf{Ax} = \mathbf{0} \Rightarrow \mathbf{A}^T \mathbf{Ax} = \mathbf{0}$ , and so  $\text{Null}(\mathbf{A}) \subset \text{Null}(\mathbf{A}^T \mathbf{A})$ . For inclusion the other way, realize that if  $\mathbf{A}^T \mathbf{Ax} = \mathbf{0}$ , then  $\mathbf{x}^T \mathbf{A}^T \mathbf{Ax} = 0$  which means  $\mathbf{Ax} = \mathbf{0}$ , since  $\mathbf{w}^T \mathbf{w} = 0$  if and only if  $\mathbf{w} = \mathbf{0}$ , and so  $\text{Null}(\mathbf{A}^T \mathbf{A}) \subset \text{Null}(\mathbf{A})$ .

**Proposition.** Let  $\mathbf{A}$  be an  $M \times N$  matrix. Then  $\text{Col}(\mathbf{A}^T \mathbf{A}) = \text{Row}(\mathbf{A})$ .

**Proof.** It is clear that if two subspaces of  $\mathbb{R}^N$  have the same orthogonal complement, then they are equal to one another. The proposition follows from the facts that

1.  $\text{Null}(\mathbf{A})$  is the orthogonal complement of  $\text{Row}(\mathbf{A})$ ,
2.  $\text{Null}(\mathbf{A}^T \mathbf{A})$  is the orthogonal complement of  $\text{Row}(\mathbf{A}^T \mathbf{A}) = \text{Col}(\mathbf{A}^T \mathbf{A})$  (where this equality comes from the fact that  $\mathbf{A}^T \mathbf{A}$  is symmetric), and
3.  $\text{Null}(\mathbf{A}^T \mathbf{A}) = \text{Null}(\mathbf{A})$ .