

week12-01-least-squares

April 7, 2024

1 George McNinch Math 87 - Spring 2024

2 Week 12

3 Least squares

4 Linear Least Squares

We are going to begin our discussion of “least squares” approximation with an example.

4.1 Example

Consider a stretch of highway with four distinct reference points **A**, **B**, **C**, and **D**:

[**A**] — x_1 — [**B**] — x_2 — [**C**] — x_3 — [**D**]

Write $x_1 = \mathbf{AB}$ for the distance from **A** to **B**, $x_2 = \mathbf{BC}$, $x_3 = \mathbf{CD}$.

We take some measurements – which potentially reflect errors – , and we seek the *best approximation* to the distances x_1 , x_2 , x_3 .

The measurements taken are as follows:

segment	AD	AC	BD	AB	CD
length	89 m	67 m	53 m	35 m	20 m

Thus the observations suggest the following equations:

$$(1) \quad x_1 + x_2 + x_3 = 89$$

$$(2) \quad x_1 + x_2 = 67$$

$$(3) \quad x_2 + x_3 = 53$$

$$(4) \quad x_1 = 35$$

$$(5) \quad x_3 = 20$$

These equations aren’t compatible, though. Note e.g. that equations (3) -- (5) indicate the following:

$$x_1 = 35, \quad x_3 = 20, \quad x_2 = 53 - 20 = 33$$

but then we find that

$$x_1 + x_2 + x_3 = 35 + 33 + 20 = 88$$

which is incompatible with (1).

And we find that

$$x_1 + x_2 = 35 + 33 = 68$$

which is incompatible with (2).

Let's formulate these equalities in matrix form.

Thus let

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 89 \\ 67 \\ 53 \\ 35 \\ 20 \end{pmatrix}.$$

With these notations, the above equations suggest that $A\mathbf{x}$ should be equal to \mathbf{b} .

Our observation(s) in the preceding slides show, however, that the system of equations $A\mathbf{x} = \mathbf{b}$ is *inconsistent* (i.e. there is no vector \mathbf{x} which makes the equation true).

4.2 Residual

In general, given an $m \times n$ matrix A , a column vector $\mathbf{b} \in \mathbb{R}^m$ and an equation $A\mathbf{x} = \mathbf{b}$, we instead look at the so-called *residual*

$$\mathbf{r} = \mathbf{b} - A\mathbf{x},$$

and *minimize* this residual.

More precisely, we want to minimize the *magnitude* (or *length*) of this vector.

Thus if $\mathbf{r} = (r_1 \ \cdots \ r_m)^T$, we must minimize the quantity

$$\|\mathbf{r}\| = \left(\sum_{i=1}^m r_i^2 \right)^{1/2}$$

Here $\|\mathbf{r}\|$ is the magnitude, also called the Euclidean norm, of the vector \mathbf{r} .

In fact, because $f(x) = \sqrt{x}$ is an increasing function of x , we instead minimize the *square* of the magnitude of \mathbf{r} :

$$\|\mathbf{r}\|^2 = \sum_{i=1}^m r_i^2$$

Thus, we wish to find

$$\min_{\mathbf{x}} \|\mathbf{r}\|^2 = \min_{\mathbf{x}} \|\mathbf{b} - A\mathbf{x}\|^2 = \min_{x_1, x_2, \dots, x_n} \sum_{i=1}^m \left(b_i - \sum_{j=1}^n A_{ij} x_j \right)^2$$

The idea behind this minimization is to first compute for $1 \leq k \leq n$ the partial derivatives $\frac{\partial F}{\partial x_k}$ of the function

$$F(x_1, x_2, \dots, x_n) = \sum_{i=1}^m \left(b_i - \sum_{j=1}^n A_{ij} x_j \right)^2$$

Critical points - and thus possible minima - for F occur at points \mathbf{x} for which all $\frac{\partial F}{\partial x_k}(\mathbf{x}) = 0$.

Now,

$$\frac{\partial F}{\partial x_k} = \sum_{i=1}^m 2 \left(b_i - \sum_{j=1}^n A_{ij} x_j \right) (-A_{ik}) = 2 \left(\sum_{i=1}^m (-A_{ik} b_i) + \sum_{i=1}^m A_{ik} \sum_{j=1}^n A_{ij} x_j \right)$$

and this expression is equal to the k -th coefficient of the vector

$$2(-A^T \mathbf{b} + A^T A \mathbf{x})$$

Thus, the condition $\frac{\partial F}{\partial x_k} = 0$ for all k is equivalent to the so-called *normal equations*:

$$(\diamond) \quad A^T A \mathbf{x} = A^T \mathbf{b}.$$

Thus the solutions \mathbf{x} to the normal equations (\diamond) are precisely the critical points of the function F .

Recall that $A \in \mathbb{R}^{m \times n}$. Thus, $A^T \in \mathbb{R}^{n \times m}$ so that the matrix $A^T A$ is $n \times n$; in particular, $A^T A$ is always a square matrix.

Moreover, $A^T A$ is *symmetric*, since

$$(A^T A)^T = A^T (A^T)^T = A^T A.$$

We are interested here in the case of *overdetermined systems* – i.e. in the case where A has more rows than columns (“more equations than variables”). Thus $m \geq n$.

We also are interested in the case where A has rank n – i.e. A has n linearly independent columns – since otherwise we don’t expect to have enough information to find \mathbf{x} .

4.3 Proposition

Let $A \in \mathbb{R}^{m \times n}$, suppose that $m \geq n$ and that A has rank n . Then $A^T \cdot A$ is invertible.

Proof:

Since $A^T \cdot A$ is an $n \times n$ square matrix, the proposition will follow if we argue that the null space $\text{Null}(A^T A)$ is zero. So: suppose that $\mathbf{v} \in \text{Null}(A^T A)$.

Thus $A^T A \mathbf{v} = 0$ and thus also $\mathbf{v}^T A^T \cdot A \mathbf{v} = 0$.

Now,

$$0 = \mathbf{v}^T A^T \cdot A \mathbf{v} = (A \mathbf{v})^T (A \mathbf{v})$$

and of course for any vector \mathbf{w} , we know that

$$0 = \mathbf{w}^T \mathbf{w} \implies \mathbf{w} = \mathbf{0}.$$

We now conclude that $A\mathbf{v} = 0$, so $\mathbf{v} \in \text{Null}(A)$. Since A has rank n , the Null space of A is equal to zero, and we conclude that $\mathbf{v} = \mathbf{0}$.

We have now proved that $\text{Null}(A^T A)$ is zero, as required.

Remark: What we have really showed is that the symmetric matrix $A^T A$ is *definite*: $\mathbf{v}^T A^T A \mathbf{v} = 0 \implies \mathbf{v} = \mathbf{0}$.

We finally claim that this solution must minimize the magnitude of the residual.

This depends on a “second derivative test” argument for which I’m not going to give full details. The main point is that the “second derivative” in this context – known as the Hessian – coincides with the matrix $2A^T A$. Now, under our assumptions the matrix $A^T A$ is positive definite, and it follows that \mathbf{x}_0 is a global minimum for the magnitude of the residual!

4.4 Return to the example

Recall that

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 89 \\ 67 \\ 53 \\ 35 \\ 20 \end{pmatrix}.$$

So to minimize the magnitude of the residual, we must solve the normal equations:

$$A^T A \mathbf{x} = A^T \mathbf{b}.$$

Now

$$A^T A = \begin{pmatrix} 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 3 & 2 & 1 \\ 2 & 3 & 2 \\ 1 & 2 & 3 \end{pmatrix}$$

and

$$A^T \mathbf{b} = \begin{pmatrix} 191 \\ 209 \\ 162 \end{pmatrix}$$

So we need to solve the equation

$$\begin{pmatrix} 3 & 2 & 1 \\ 2 & 3 & 2 \\ 1 & 2 & 3 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 191 \\ 209 \\ 162 \end{pmatrix}$$

```
[1]: import numpy as np
import numpy.linalg as la

A= np.array([[1,1,1],[1,1,0],[0,1,1],[1,0,0],[0,0,1]])
b = np.array([89,67,53,35,20])

x0=la.solve(A.T @ A, A.T @ b)

def residual(x):
    return b - A @ x

def magnitude(x):
    return np.sqrt(x@x)

[x0,magnitude(residual(x0))]
```

```
[1]: [array([35.125, 32.5  , 20.625]), 1.1726039399558574]
```

Thus the *least squares solution* is

$$\mathbf{x}_0 = \begin{pmatrix} 35.125 \\ 32.5 \\ 20.625 \end{pmatrix}$$

and

$$\|\mathbf{b} - A\mathbf{x}_0\| \approx 1.1726$$

Recall that our “first guess” for a solution (based on some of the measurements) was

$$\mathbf{x}_1 = \begin{pmatrix} 35 \\ 33 \\ 20 \end{pmatrix}$$

The residual is indeed larger for x_1 :

```
[2]: x1 = np.array([35,33,20])

magnitude(residual(x1))
```

```
[2]: 1.4142135623730951
```

Let’s note that `numpy` already implements this least-squares functionality:

– you can [read more about it here](#)

```
[3]: res=la.lstsq(A,b,rcond=None)
      res[0]
```

```
[3]: array([35.125, 32.5  , 20.625])
```

5 Example recapitulated

Consider a group of 10 employees $[a,b,c,\dots,h,i,j]$ of a certain company.

```
[ ]:
```