

Least Squares and Projections

February 26, 2021

Least Squares

We want to focus on solving $A\mathbf{x} = \mathbf{b}$ for an over-determined system so A is $m \times n$ with $m > n$ (more equations than variables) and where $\text{rank}(A) = n$ so the columns are linearly independent. In general, there will not be a solution to this so we ask instead, what is the "best possible" solution. This of course is vague so to be less vague consider all possible values of $A\mathbf{x}$, this is just $\text{Im}(A) = \text{CS}(A)$ and we want to find the point in $\text{CS}(A)$ closest to \mathbf{b} where closest is measured in the usual notion of distance. So we mean to find $\hat{\mathbf{x}}$ so that $\|A\hat{\mathbf{x}} - \mathbf{b}\|_2$ is as small as possible. This is equivalent to minimizing $\|A\hat{\mathbf{x}} - \mathbf{b}\|_2^2 = \langle A\hat{\mathbf{x}} - \mathbf{b}, A\hat{\mathbf{x}} - \mathbf{b} \rangle$, this is where the name "least squares" comes from.

Consider

$$\begin{aligned} \|A\mathbf{x} - \mathbf{b}\|_2^2 &= \langle A\mathbf{x} - \mathbf{b}, A\mathbf{x} - \mathbf{b} \rangle \\ &= \end{aligned}$$

Think geometrically, we want to find $\hat{\mathbf{b}}$ in $\text{CS}(A)$ that is as close as possible to \mathbf{b} . It would make sense that $\hat{\mathbf{b}}$ would be the *orthogonal projection* of \mathbf{b} onto $\text{CS}(A)$. That is find $\hat{\mathbf{b}} \in \text{CS}(A)$ so that $\mathbf{b} - \hat{\mathbf{b}}$ is orthogonal to $\text{CS}(A)$.

Notice $\mathbf{b} - \hat{\mathbf{b}} \perp \text{CS}(A)$ iff $A^T(\mathbf{b} - \hat{\mathbf{b}}) = \mathbf{0}$ so we want to solve

$$\begin{aligned} A^T\mathbf{b} &= A^T\hat{\mathbf{b}} && \text{since } \mathbf{b} - \hat{\mathbf{b}} \perp \text{CS}(A) \\ \hat{\mathbf{b}} &= A\mathbf{u} \text{ for some } \mathbf{u} \in \mathbb{R}^n && \text{since } \hat{\mathbf{b}} \in \text{CS}(A) \end{aligned}$$

This leads to

$$A^T\mathbf{b} = A^T\hat{\mathbf{b}} = A^TA\mathbf{u}$$

so

$$\mathbf{u} = (A^TA)^{-1}A^T\mathbf{b}$$

Exercise 1. Show that $A^T A$ is non-singular. Recall A is $m \times n$, $n < m$, and $\text{rank}(A) = n$.

In this way we have solved the orthogonal projection problem and the least squares problem in one go. The least squares solution to $A\mathbf{x} = \mathbf{b}$ is

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

and the orthogonal projection of \mathbf{b} onto $\text{CS}(A)$ is

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

The $m \times m$ matrix

$$P = A(A^T A)^{-1} A^T$$

is called the *orthogonal projection matrix* of \mathbb{R}^m onto $\text{CS}(A)$ and $P\mathbf{b}$ is the orthogonal projection of \mathbf{b} onto $\text{CS}(A)$ for any $\mathbf{b} \in \mathbb{R}^m$.

Fitting polynomials to data.

Given a bunch of data in \mathbb{R}^2 of the form (x_i, i_i) for $i = 1, \dots, N$ we can try to fit a polynomial of order m (usually much smaller than N) to the data as follows. We'd like to find $\alpha_0, \alpha_1, \dots, \alpha_m$ so that $y_i = \alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 + \dots + \alpha_m x_i^m$ for all $i = 1, \dots, N$. This can be written in matrix form as: Find a vector $\boldsymbol{\alpha}$ so that

$$\begin{bmatrix} 1 & \mathbf{x} & \mathbf{x}^2 & \dots & \mathbf{x}^m \end{bmatrix} \boldsymbol{\alpha} = \mathbf{y}$$

Here $\mathbf{y} = (y_1, \dots, y_N)^T$ and $\mathbf{x}^k = (x_1^k, \dots, x_N^k)^T$. Typically if the X_i 's are somewhat random (maybe even distinct) and $m < N$, the coefficient matrix will have rank $m + 1$ and thus we can apply the technique from above to find the least squares solution to this. (Typically there will not be an actual solution!)

Let $A = \begin{bmatrix} 1 & \mathbf{x} & \mathbf{x}^2 & \dots & \mathbf{x}^m \end{bmatrix}$ and $\hat{\boldsymbol{\alpha}} = (A^T A)^{-1} A^T \mathbf{y}$, then the m^{th} degree polynomial that best fits the data is $\hat{\boldsymbol{\alpha}}^T \mathbf{x} = \sum_{i=1}^N \hat{\alpha}_i x^i$ (here $\mathbf{x} = (1, x, x^2, x^3, \dots, x^m)$).

The case of best fitting line is just where $m = 1$.

This is super easy to implement in Octave/MATLAB!

```

1 N = 600;
2 M = 17;
3
4 % Generate N uniformly distributed x values
5 % between -4 and 4.
6 x = 8*rand(N,1) - 4;
7
8 X = sort(x);
9
10 % Apply some function to the x values
11 Y = sin(4*X) + cos(3*X) - X/8; % green

```

```

12 % Add some noise (our simulated data)
13 y = Y + 2*rand(N,1) -1; % blue
14
15 % Build the matrix [1 x x^2 ... x^M]
16 A = zeros(N,M);
17
18 for k = 0:M
19     A(:,k+1) = X.^k;
20 end
21
22 % find the coefficients of our M-degree poly
23 alpha = (A'*A)^-1*A'*y;
24
25 % Generate values based on our polynomial (red)
26 haty = A*alpha;
27
28 plot(X,y,"b.",X,haty,'r-',X,Y,'g-')
29
30 axis([-4 4 -1 1])
31 axis('square')

```

QR decomposition and least squares.

Any $m \times n$ matrix A of rank n where $n < m$, can be written as QR where Q is orthogonal $m \times n$ and R is upper triangular $n \times n$ (invertible). Recall when finding the least square solution to $A\mathbf{x} = \mathbf{b}$ we had

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

This is the same as solving

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

which is equivalent to

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

and this reduces to

$$R \mathbf{x} = Q^T \mathbf{b}$$

by multiplying both sides by $(R^T)^{-1}$ which is trivial to solve by back substitution, since R is upper triangular.

Getting the QR decomposition really just follows from the Gram-Schmidt procedure applied to the columns of A (which are assumed to be linearly independent.) Recall in GS we simply subtract the orthogonal projection of \mathbf{a}_i onto $\text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_{i-1}\}$ from \mathbf{a}_i itself, that is:

$$\mathbf{q}_i = \mathbf{a}_i - A_{i-1}(A_{i-1}A_{i-1}^T)^{-1}A_{i-1}^T \mathbf{a}_i,$$

where $A_j = [\mathbf{a}_1 \ \dots \ \mathbf{a}_j]$

To make the \mathbf{q}_i 's unit vectors simply normalize them setting $\hat{\mathbf{q}}_i = \mathbf{q}_i / \|\mathbf{q}_i\|$.

Recall, that $\text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_j\} = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_j\}$ by construction so

$$\mathbf{a}_i = \langle \hat{\mathbf{q}}_1, \mathbf{a}_i \rangle \hat{\mathbf{q}}_1 + \dots + \langle \hat{\mathbf{q}}_i, \mathbf{a}_i \rangle \hat{\mathbf{q}}_i.$$

This clearly shows that $A = QR$ where $Q = [\hat{\mathbf{q}}_1 \ \dots \ \hat{\mathbf{q}}_n]$ and hence $Q^T R = Q^T QR = IR = R$, since Q is unitary, that is $R = Q^T A$.

This yields very simple MATLAB code:

```

1 function [Q,R] = QR(A)
2
3 % Usage: [Q,R] = QR(A)
4 % Assumption: A is a rank m, m x n matrix
5 % Returns: [Q,R], Q is unitary m x n, R is upper triangular n x n
6
7 [m, n] = size(A);
8
9 Q = zeros(m,n);
10 R = zeros(n,n);
11
12 % Just normalize the first vector
13 q = A(:,1);
14 q = q/(q'*q)^.5;
15 Q(:,1) = q;
16
17 % Run GS
18 for i = 2:n
19     B = A(:,1:i-1);
20     q = A(:,i);
21     q = q - B*(B'*B)^-1*B'*q;
22     q = q/(q'*q)^.5;
23     Q(:,i) = q;
24 end
25
26 R = Q'*A
27
28 end

```