

LEAST SQUARES APPROXIMATION

SYNOPSIS. Having established the key theoretical results, the main linear algebra problems can now be solved. The typical scenario within data science applications is that a large dimensional vector representation of some object is available and greater insight is sought by seeking a description in terms of linear combination of a small number of vectors. The large dimensional object might not be exactly recovered and the focus is on obtaining the best possible approximation. In Euclidean spaces with distances measured by the 2-norm, best approximants are readily found by the generalization of the Pythagorean theorem to high-dimensional spaces.

1. Orthogonal projection

Consider a partition of a vector space U into orthogonal subspaces $U = V \oplus W$, $V = W^\perp$, $W = V^\perp$, typically $U = \mathbb{R}^m$, $V \subset \mathbb{R}^m$, $W \subset \mathbb{R}^m$, $\dim V = n$, $\dim W = m - n$. If $\mathbf{V} = [\mathbf{v}_1 \dots \mathbf{v}_n] \in \mathbb{R}^{m \times n}$ is a basis for V and $\mathbf{W} = [\mathbf{w}_1 \dots \mathbf{w}_{m-n}] \in \mathbb{R}^{m \times (m-n)}$ is a basis for W , then $\mathbf{U} = [\mathbf{v}_1 \dots \mathbf{v}_n \mathbf{w}_1 \dots \mathbf{w}_{m-n}]$ is a basis for U . Even though the matrices \mathbf{V} , \mathbf{W} are not necessarily square, they are said to be orthonormal when all columns are of unit norm and orthogonal to one another. In this case computation of the matrix product $\mathbf{V}^T \mathbf{V}$ leads to the formation of the identity matrix within \mathbb{R}^n

$$\mathbf{V}^T \mathbf{V} = \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_n] = \begin{bmatrix} \mathbf{v}_1^T \mathbf{v}_1 & \mathbf{v}_1^T \mathbf{v}_2 & \dots & \mathbf{v}_1^T \mathbf{v}_n \\ \mathbf{v}_2^T \mathbf{v}_1 & \mathbf{v}_2^T \mathbf{v}_2 & \dots & \mathbf{v}_2^T \mathbf{v}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{v}_n^T \mathbf{v}_1 & \mathbf{v}_n^T \mathbf{v}_2 & \dots & \mathbf{v}_n^T \mathbf{v}_n \end{bmatrix} = \mathbf{I}_n.$$

Similarly, $\mathbf{W}^T \mathbf{W} = \mathbf{I}_{m-n}$. Whereas for the square orthogonal matrix \mathbf{U} multiplication both on the left and the right by its transpose leads to the formation of the identity matrix

$$\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_m,$$

the same operations applied to rectangular orthogonal matrices lead to different results

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}_n, \mathbf{V} \mathbf{V}^T = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_n] \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} = \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T, \text{rank}(\mathbf{v}_i \mathbf{v}_i^T) = 1$$

A simple example is provided by taking $\mathbf{V} = \mathbf{I}_{m,n}$, the first n columns of the identity matrix in which case

$$\mathbf{V} \mathbf{V}^T = \sum_{i=1}^n \mathbf{e}_i \mathbf{e}_i^T = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{m \times m}.$$

Applying $\mathbf{P} = \mathbf{V} \mathbf{V}^T$ to some vector $\mathbf{b} \in \mathbb{R}^m$ leads to a vector $\mathbf{r} = \mathbf{P} \mathbf{b}$ whose first n components are those of \mathbf{b} , and the remaining $m - n$ are zero. The subtraction $\mathbf{b} - \mathbf{r}$ leads to a new vector $\mathbf{s} = (\mathbf{I} - \mathbf{P}) \mathbf{b}$ that has the first components equal to zero, and the remaining $m - n$ the same as those of \mathbf{b} . Such operations are referred to as *projections*, and for $\mathbf{V} = \mathbf{I}_{m,n}$ correspond to projection onto the span $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$.

Returning to the general case, the orthogonal matrices $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{m \times n}$, $W \in \mathbb{R}^{m \times (m-n)}$ are associated with linear mappings $b = f(x) = Ux$, $r = g(b) = Pb$, $s = h(b) = (I - P)b$. The mapping f gives the components in the I basis of a vector whose components in the U basis are x . The mappings g, h project a vector onto $\text{span}\{v_1, \dots, v_n\}$, $\text{span}\{w_1, \dots, w_{m-n}\}$, respectively. When V, W are orthogonal matrices the projections are also orthogonal $r \perp s$. Projection can also be carried out onto nonorthogonal spanning sets, but the process is fraught with possible error, especially when the angle between basis vectors is small, and will be avoided henceforth. Notice that projection of a vector already in the spanning set simply returns the same vector, which leads to a general definition.

DEFINITION. The mapping is called a **projection** if $f \circ f = f$, or if for any $u \in U$, $f(f(u)) = f(u)$. With P the matrix associated f , a projection matrix satisfies $P^2 = P$.

Orthogonal projections onto the column space $C(Q)$ of an orthonormal matrix are of great practical utility, and satisfy the above definition

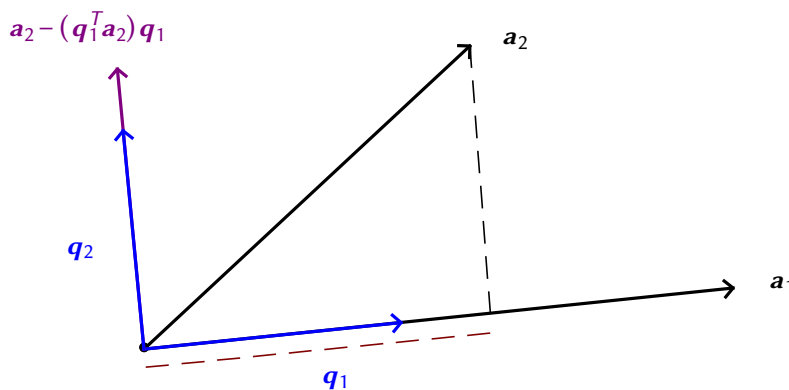
$$P_Q = QQ^T$$

$$P_Q^2 = P_Q P_Q = QQ^T QQ^T = Q(Q^T Q)Q^T = QIQ^T = QQ^T = P_Q.$$

2. Gram-Schmidt algorithm

Orthonormal vector sets $\{q_1, \dots, q_n\}$ are of the greatest practical utility, leading to the question of whether some such a set can be obtained from an arbitrary set of vectors $\{a_1, \dots, a_n\}$. This is possible for independent vectors, through what is known as the Gram-Schmidt algorithm

1. Start with an arbitrary direction a_1
2. Divide by its norm to obtain a unit-norm vector $q_1 = a_1 / \|a_1\|$
3. Choose another direction a_2
4. Subtract off its component along previous direction(s) $a_2 - (q_1^T a_2) q_1$
5. Divide by norm $q_2 = (a_2 - (q_1^T a_2) q_1) / \|a_2 - (q_1^T a_2) q_1\|$
6. Repeat the above



$$P_1 a_2 = (q_1 q_1^T) a_2 = q_1 (q_1^T a_2) = (q_1^T a_2) q_1$$

The above geometrical description can be expressed in terms of matrix operations as

$$A = (\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n) = (\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_n) \begin{pmatrix} r_{11} & r_{12} & r_{13} & \dots & r_{1n} \\ 0 & r_{22} & r_{23} & \dots & r_{2n} \\ 0 & 0 & r_{33} & \dots & r_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \dots & r_{mn} \end{pmatrix} = QR,$$

equivalent to the system

$$\begin{cases} \mathbf{a}_1 = r_{11}\mathbf{q}_1 \\ \mathbf{a}_2 = r_{12}\mathbf{q}_1 + r_{22}\mathbf{q}_2 \\ \vdots \\ \mathbf{a}_n = r_{1n}\mathbf{q}_1 + r_{2n}\mathbf{q}_2 + \dots + r_{nn}\mathbf{q}_n \end{cases}$$

The system is easily solved by *forward substitution* resulting in what is known as the (modified) *Gram-Schmidt algorithm*, transcribed below both in pseudo-code and in Julia.

Algorithm (Gram-Schmidt)

Given n vectors $\mathbf{a}_1, \dots, \mathbf{a}_n$
 Initialize $\mathbf{q}_1 = \mathbf{a}_1, \dots, \mathbf{q}_n = \mathbf{a}_n, \mathbf{R} = \mathbf{I}_n$
 for $i = 1$ to n
 $r_{ii} = (\mathbf{q}_i^T \mathbf{q}_i)^{1/2}$
 if $r_{ii} < \epsilon$ break;
 $\mathbf{q}_i = \mathbf{q}_i / r_{ii}$
 for $j = i+1$ to n
 $r_{ij} = \mathbf{q}_i^T \mathbf{a}_j; \mathbf{q}_j = \mathbf{q}_j - r_{ij}\mathbf{q}_i$
 end
 end
 return \mathbf{Q}, \mathbf{R}

```
∴ function mgs(A)
    m,n=size(A); Q=A; R=eye(n);
    for i=1:n
        R[i,i] = sqrt(Q[:,i]'*Q[:,i]);
        if (R[i,i]<eps) break;
        Q[:,i] = Q[:,i]/R[i,i];
        for j=i+1:n
            R[i,j] = Q[:,i]'*A[:,j];
            Q[:,j] = Q[:,j] - R[i,j]*Q[:,i];
        end;
    end;
    return Q,R
end
```

The input matrix A might have linearly dependent columns in which case $r_{ii} \approx 0$ for some i , and the if-instruction interrupts the algorithm. The Gram-Schmidt algorithm furnishes a *factorization*

$$QR = A$$

with $\mathbf{Q} \in \mathbb{R}^{m \times n}$ an orthonormal matrix and $\mathbf{R} \in \mathbb{R}^{n \times n}$ an upper triangular matrix, known as the *QR-factorization*. Since the column vectors within \mathbf{Q} were obtained through linear combinations of the column vectors of A , $C(A) = C(Q)$. The *QR-factorization* is of great utility in solving problems within linear algebra.

3. Least squares problems in \mathbb{R}^m

3.1. Problem formulation and solution by orthogonal projection

A typical situation in applications is that a vector $\mathbf{y} \in \mathbb{R}^m$ represents a complex object with $m \gg 1$. A simpler representation of the object is sought through a linear combination $\mathbf{v} = A\mathbf{c}$, with $A \in \mathbb{R}^{m \times n}$ and $n < m$, usually $m \ll n$ (Fig. 1).

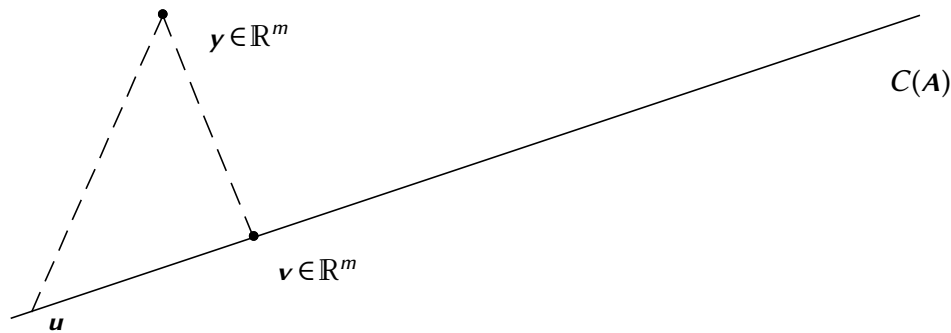


Figure 1. Least squares problem: find $\mathbf{v} \in C(\mathbf{A})$, $\mathbf{A} \in \mathbb{R}^{m \times n}$ closest to some given \mathbf{y} in the 2-norm

The magnitude of the difference between the exact object \mathbf{y} and the approximation $\mathbf{A}\mathbf{c}$ is measured through a norm, and in least squares the 2-norm is adopted. The problem is to minimize the error $\varepsilon = \|\mathbf{e}\| = \|\mathbf{y} - \mathbf{A}\mathbf{c}\| = (\mathbf{e}^T \mathbf{e})^{1/2}$. This is stated mathematically as

$$\min_{\mathbf{c} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{c}\|,$$

and within \mathbb{R}^m the minimal (2-norm) distance is obtained when $\mathbf{y} - \mathbf{v}$ is orthogonal to $C(\mathbf{A})$. Note that if another type of norm were to be adopted, the orthogonality condition would no longer necessarily hold. For the 2-norm however, the orthogonality condition leads to a straightforward solution of the problem through orthogonal projection:

1. Find an orthonormal basis for column space of \mathbf{A} by \mathbf{QR} factorization, $\mathbf{QR} = \mathbf{A}$.
2. State that \mathbf{v} is the projection of \mathbf{y} , $\mathbf{v} = \mathbf{P}_{C(\mathbf{A})}\mathbf{y} = \mathbf{P}_Q \mathbf{y} = \mathbf{Q}\mathbf{Q}^T \mathbf{y}$.
3. State that \mathbf{v} is within the column space of \mathbf{A} , $\mathbf{v} = \mathbf{A}\mathbf{c} = \mathbf{Q}\mathbf{R}\mathbf{c}$.
4. Set equal the two expressions of \mathbf{v} , $\mathbf{Q}\mathbf{Q}^T \mathbf{y} = \mathbf{Q}\mathbf{R}\mathbf{c}$. This is an equality between two linear combinations of the columns of \mathbf{Q} . For \mathbf{Q} orthonormal the scaling coefficients of the two linear combinations must be equal leading to $\mathbf{R}\mathbf{c} = \mathbf{Q}^T \mathbf{y}$.
5. Solve the triangular system to find \mathbf{c} .

3.2. Linear regression

The approach to compressing data $D = \{(x_i, y_i) \mid i = 1, \dots, m\}$ suggested by calculus concepts is to form the sum of squared differences between $y(x_i)$ and y_i , for example for $y(x) = c_0 + c_1 x$ when carrying out linear regression,

$$S(c_0, c_1) = \sum_{i=1}^m (y(x_i) - y_i)^2 = \sum_{i=1}^m (c_0 + c_1 x_i - y_i)^2$$

and seek (c_0, c_1) that minimize $S(c_0, c_1)$. The function $S(c_0, c_1) \geq 0$ can be thought of as the height of a surface above the $c_0 c_1$ plane, and the gradient ∇S is defined as a vector in the direction of steepest slope. When at some point on the surface if the gradient is different from the zero vector $\nabla S \neq \mathbf{0}$, travel in the direction of the gradient would increase the height, and travel in the opposite direction would decrease the height. The minimal value of S would be attained

when no local travel could decrease the function value, which is known as stationarity condition, stated as $\nabla S = 0$. Applying this to determining the coefficients (c_0, c_1) of a linear regression leads to the equations

$$\frac{\partial S}{\partial c_0} = 0 \Rightarrow 2 \sum_{i=1}^m (c_0 + c_1 x_i - y_i) = 0 \Leftrightarrow m c_0 + \left(\sum_{i=1}^m x_i \right) c_1 = \sum_{i=1}^m y_i,$$

$$\frac{\partial S}{\partial c_1} = 0 \Rightarrow 2 \sum_{i=1}^m (c_0 + c_1 x_i - y_i) x_i = 0 \Leftrightarrow \left(\sum_{i=1}^m x_i \right) c_0 + \left(\sum_{i=1}^m x_i^2 \right) c_1 = \sum_{i=1}^m x_i y_i.$$

The above calculations can become tedious, and do not illuminate the geometrical essence of the calculation, which can be brought out by reformulation in terms of a matrix-vector product that highlights the particular linear combination that is sought in a linear regression. Form a vector of errors with components $e_i = y(x_i) - y_i$, which for linear regression is $y(x) = c_0 + c_1 x$. Recognize that $y(x_i)$ is a linear combination of 1 and x_i with coefficients c_0, c_1 , or in vector form

$$\mathbf{e} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \end{bmatrix} - \mathbf{y} = [\mathbf{1} \ \mathbf{x}] \mathbf{c} - \mathbf{y} = \mathbf{A} \mathbf{c} - \mathbf{y}.$$

The norm of the error vector $\|\mathbf{e}\|$ is smallest when $\mathbf{A} \mathbf{a}$ is as close as possible to \mathbf{y} . Since $\mathbf{A} \mathbf{a}$ is within the column space of $C(\mathbf{A})$, $\mathbf{A} \mathbf{a} \in C(\mathbf{A})$, the required condition is for \mathbf{e} to be orthogonal to the column space

$$\mathbf{e} \perp C(\mathbf{A}) \Rightarrow \mathbf{A}^T \mathbf{e} = \begin{bmatrix} \mathbf{1}^T \\ \mathbf{x}^T \end{bmatrix} \mathbf{e} = \begin{bmatrix} \mathbf{1}^T \mathbf{e} \\ \mathbf{x}^T \mathbf{e} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \mathbf{0}$$

$$\mathbf{A}^T \mathbf{e} = \mathbf{0} \Leftrightarrow \mathbf{A}^T (\mathbf{A} \mathbf{c} - \mathbf{y}) = \mathbf{0} \Leftrightarrow (\mathbf{A}^T \mathbf{A}) \mathbf{c} = \mathbf{A}^T \mathbf{y} = \mathbf{b} \Leftrightarrow \mathbf{N} \mathbf{c} = \mathbf{b}.$$

The above is known as the normal system, with $\mathbf{N} = \mathbf{A}^T \mathbf{A}$ is the normal matrix. The system $\mathbf{N} \mathbf{c} = \mathbf{b}$ can be interpreted as seeking the coordinates in the $\mathbf{N} = \mathbf{A}^T \mathbf{A}$ basis of the vector $\mathbf{b} = \mathbf{A}^T \mathbf{y}$. An example can be constructed by randomly perturbing a known function $y(x) = a_0 + a_1 x$ to simulate measurement noise and compare to the approximate \tilde{c} obtained by solving the normal system.

```
∴ m=100; x=(0:m-1)./m; c0=2; c1=3; yex=c0.+c1*x; y=(yex.+rand(m,1).-0.5);
```

```
∴ A=ones(m,2); A[:,2]=x[:]; At=transpose(A); N=At*A; b=At*y;
```

```
∴ c = N\b
```

$$\begin{bmatrix} 1.9524719146017488 \\ 3.114552636517776 \end{bmatrix} \quad (1)$$

```
∴
```

3.3. Least squares polynomial approximations of data

Forming the normal system of equations can lead to numerical difficulties, especially when the columns of \mathbf{A} are close to linear dependence. It is preferable to adopt the general procedure of solving a least squares problem by projection, in which case the above linear regression becomes:

$$\mathbf{Q} \mathbf{R} = \mathbf{A}, \mathbf{R} \mathbf{c} = \mathbf{Q}^T \mathbf{y}.$$

```
∴ QR=qr(A); Q=QR.Q[:,1:2]; R=QR.R[1:2,1:2];
```

```
∴ c = R\ (transpose(Q)*y)
```

$$\begin{bmatrix} 1.9524719146017504 \\ 3.114552636517771 \end{bmatrix} \quad (2)$$

```
∴
```

The above procedure can be easily extended to define quadratic or cubic regression, the problem of finding the best polynomials of degree 2 or 3 that fit the data. Quadratic regression is simply accomplished by adding a column to A containing the squares of the x vector

$$A = [a_1 \ a_2 \ a_3] = [\mathbf{1} \ x \ x^2]$$

with the column vector $a_k = x^{k-1} \in \mathbb{R}^m$ has components x_i^{k-1} for $i = 1, 2, \dots, m$.

```
∴ m=100; x=(0:m-1)./m; c0=2; c1=3; c2=-5; yex=c0.+c1*x.+c2*x.^2;
```

```
∴ y=(yex.+rand(m,1).-0.5);
```

```
∴ A=ones(m,3); A[:,2]=x[:]; A[:,3]=x[:].^2; QR=qr(A); Q=QR.Q[:,1:3];  
R=QR.R[1:3,1:3];
```

```
∴ c = R\ (transpose(Q)*y)
```

$$\begin{bmatrix} 1.968925863819198 \\ 2.9292845098793223 \\ -4.836797382728068 \end{bmatrix} \quad (3)$$

```
∴
```