

II.2 Least Squares : Four Ways

Many applications lead to **unsolvable linear equations** $Ax = b$. It is ironic that this is such an important problem in linear algebra. We can't throw equations away, we need to produce a best solution \hat{x} . **The least squares method chooses \hat{x} to make $\|b - A\hat{x}\|^2$ as small as possible.** Minimizing that error means that its derivatives are zero: those are the *normal equations* $A^T A \hat{x} = A^T b$. Their geometry will be in Figure II.2.

This section explains four ways to solve those important (and solvable!) equations:

- 1 The SVD of A leads to its **pseudoinverse** A^+ . Then $\hat{x} = A^+ b$: One short formula.
- 2 $A^T A \hat{x} = A^T b$ can be solved directly when A has **independent columns**.
- 3 The Gram-Schmidt idea produces **orthogonal columns** in Q . Then $A = QR$.
- 4 **Minimize** $\|b - Ax\|^2 + \delta^2 \|x\|^2$. That penalty changes the normal equations to $(A^T A + \delta^2 I)x_\delta = A^T b$. Now the matrix is invertible and x_δ goes to \hat{x} as $\delta \rightarrow 0$.

$A^T A$ has an attractive symmetry. But its size may be a problem. And its condition number—measuring the danger of unacceptable roundoff error—is the *square* of the condition number of A . In well-posed problems of moderate size we go ahead to solve the least squares equation $A^T A \hat{x} = A^T b$, but in large or ill-posed problems we find another way.

We could orthogonalize the columns of A . We could use its SVD. For really large problems we sample the column space of A by simply multiplying Av for **random vectors** v . This seems to be the future for very big computations: a high probability of success.

First of all, please allow us to emphasize the importance of $A^T A$ and $A^T C A$. That matrix C is often a positive diagonal matrix. It gives stiffnesses or conductances or edge capacities or inverse variances $1/\sigma^2$ —the constants from science or engineering or statistics that define our particular problem: the “weights” in weighted least squares.

Here is a sample of the appearances of $A^T A$ and $A^T C A$ in applied mathematics:

In mechanical engineering, $A^T A$ (or $A^T C A$) is the **stiffness matrix**

In circuit theory, $A^T A$ (or $A^T C A$) is the **conductance matrix**

In graph theory, $A^T A$ (or $A^T C A$) is the (weighted) **graph Laplacian**

In mathematics, $A^T A$ is the **Gram matrix**: inner products of columns of A

In large problems, $A^T A$ is expensive and often dangerous to compute. We avoid it if we can! The Gram-Schmidt way replaces A by QR (orthogonal Q , triangular R). Then $A^T A$ is the same as $R^T Q^T Q R = R^T R$. And the fundamental equation $A^T A \hat{x} = A^T b$ becomes $R^T R \hat{x} = R^T Q^T b$. Finally this is $R \hat{x} = Q^T b$, safe to solve and fast too.

Thus $A^T A$ and $A^T C A$ are crucial matrices—but paradoxically, we try not to compute them. Orthogonal matrices and triangular matrices: Those are the good ones.

A^+ is the Pseudoinverse of A

I will first describe the pseudoinverse A^+ in words. If A is invertible then A^+ is A^{-1} . If A is m by n then A^+ is n by m . When A multiplies a vector x in its row space, this produces Ax in the column space. Those two spaces have equal dimension r (the rank). Restricted to these spaces A is always invertible—and A^+ inverts A . Thus $A^+Ax = x$ exactly when x is in the row space. And $AA^+b = b$ when b is in the column space.

The nullspace of A^+ is the nullspace of A^T . It contains the vectors y in \mathbf{R}^m with $A^T y = 0$. Those vectors y are perpendicular to every Ax in the column space. For these y , we accept $x^+ = A^+y = 0$ as the best solution to the unsolvable equation $Ax = y$. Altogether A^+ inverts A where that is possible :

The pseudoinverse of $A = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$ is $A^+ = \begin{bmatrix} 1/2 & 0 \\ 0 & 0 \end{bmatrix}$.

The whole point is to produce a suitable “pseudoinverse” when A has no inverse.

Rule 1 If A has independent columns, then $A^+ = (A^T A)^{-1} A^T$ and so $A^+ A = I$.

Rule 2 If A has independent rows, then $A^+ = A^T (A A^T)^{-1}$ and so $AA^+ = I$.

Rule 3 A diagonal matrix Σ is inverted where possible—otherwise Σ^+ has zeros :

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \Sigma^+ = \begin{bmatrix} 1/\sigma_1 & 0 & 0 \\ 0 & 1/\sigma_2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad \begin{array}{l} \text{On the four subspaces} \\ \Sigma^+ \Sigma = I \quad \Sigma \Sigma^+ = I \\ \Sigma^+ \Sigma = 0 \quad \Sigma \Sigma^+ = 0 \end{array}$$

All matrices

The pseudoinverse of $A = U \Sigma V^T$ is $A^+ = V \Sigma^+ U^T$. (1)

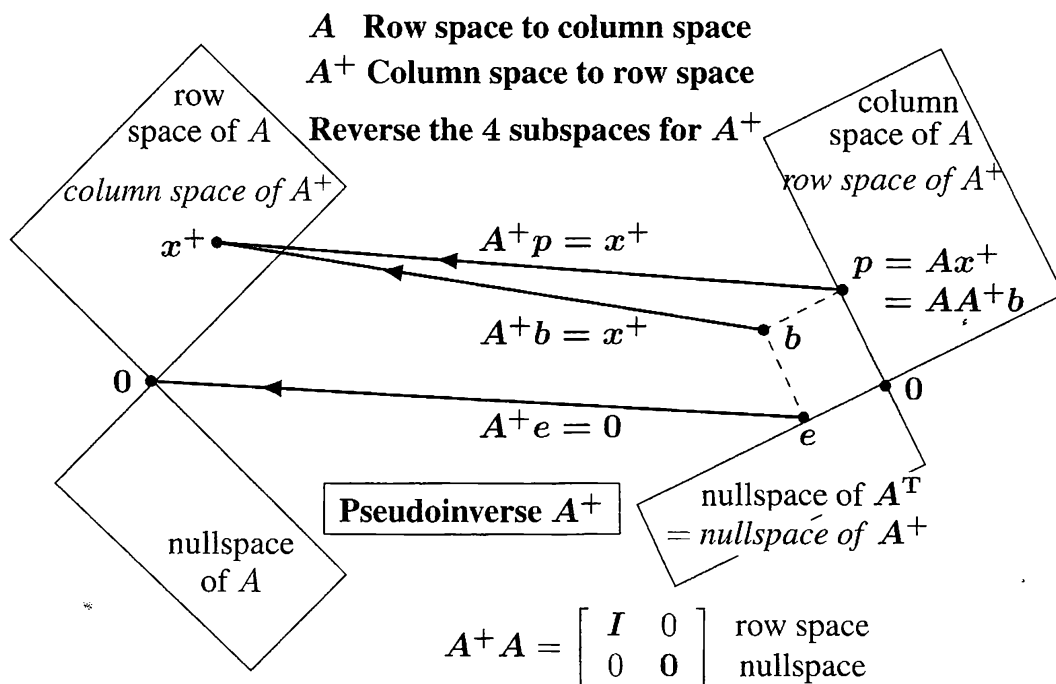


Figure II.1: Vectors $p = Ax^+$ in the column space of A go back to x^+ in the row space.

This pseudoinverse A^+ (sometimes written A^\dagger with a dagger instead of a plus sign) solves the least squares equation $A^T A \hat{x} = A^T b$ in one step. This page verifies that $x^+ = A^+ b = V \Sigma^+ U^T b$ is best possible. At the end of this section, we look in more detail at A^+ .

Question: The formula $A^+ = V \Sigma^+ U^T$ uses the SVD. Is the SVD essential to find A^+ ?

Answer: No, A^+ could also be computed directly from A by modifying the elimination steps that usually produce A^{-1} . However each step of arithmetic would have to be exact! You need to distinguish exact zeros from small nonzeros. That is the hard part of A^+ .

The Least Squares Solution to $Ax = b$ is $x^+ = A^+ b$

I have written x^+ instead of \hat{x} because the vector x^+ has two properties:

- 1 $x = x^+ = A^+ b$ makes $\|b - Ax\|^2$ as small as possible. **Least squares solution**
- 2 If another \hat{x} achieves that minimum then $\|x^+\| < \|\hat{x}\|$. **Minimum norm solution**

$x^+ = A^+ b$ is the **minimum norm least squares solution**. When A has independent columns and rank $r = n$, this is the only least squares solution. But if there are nonzero vectors x in the nullspace of A (so $r < n$), they can be added to x^+ . The error $b - A(x^+ + x)$ is not affected when $Ax = 0$. But the length $\|x^+ + x\|^2$ will grow to $\|x^+\|^2 + \|x\|^2$. Those pieces are orthogonal: Row space \perp nullspace.

So the minimum norm (shortest) solution of $A^T A \hat{x} = A^T b$ is $x^+ = A^+ b$, x^+ has a zero component in the nullspace of A .

Example 1 The shortest least squares solution to $\begin{bmatrix} 3 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 8 \end{bmatrix}$ is x^+
 $x^+ = A^+ b = \begin{bmatrix} 1/3 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 6 \\ 8 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$. All vectors $\begin{bmatrix} 0 \\ x_2 \end{bmatrix}$ are in the nullspace of A .

All the vectors $\hat{x} = \begin{bmatrix} 2 \\ x_2 \end{bmatrix}$ minimize $\|b - A\hat{x}\|^2 = 64$. But $x^+ = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ is shortest.

That example shows the least squares solutions when A is a diagonal matrix like Σ . To allow every matrix $U \Sigma V^T$, we have to account for the orthogonal matrices U and V . We can freely multiply by U^T without changing any lengths, because $U^T U = I$:

$$\text{Squared error} \quad \|b - Ax\|^2 = \|b - U \Sigma V^T x\|^2 = \|U^T b - \Sigma V^T x\|^2. \quad (2)$$

Set $w = V^T x$ to get $\|U^T b - \Sigma w\|^2$. **The best w is $\Sigma^+ U^T b$.** And finally x^+ is $A^+ b$:

$$w = V^T x^+ = \Sigma^+ U^T b \quad \text{and} \quad V^T = V^{-1} \quad \text{lead to} \quad x^+ = V \Sigma^+ U^T b = A^+ b. \quad (3)$$

The SVD solved the least squares problem in one step $A^+ b$. The only question is the computational cost. Singular values and singular vectors cost more than elimination. The next two proposed solutions work directly with the linear equations $A^T A \hat{x} = A^T b$. This succeeds when $A^T A$ is invertible—and then \hat{x} is the same as x^+ .

When is $A^T A$ Invertible ?

The invertibility (or not) of the matrix $A^T A$ is an important question with a nice answer : **$A^T A$ is invertible exactly when A has independent columns. If $Ax = 0$ then $x = 0$**

Always A and $A^T A$ have the same nullspace ! This is because $A^T Ax = 0$ always leads to $x^T A^T Ax = 0$. This is $\|Ax\|^2 = 0$. Then $Ax = 0$ and x is in $N(A)$. For all matrices :

$$N(A^T A) = N(A) \text{ and } C(AA^T) = C(A) \text{ and } \mathbf{rank}(A^T A) = \mathbf{rank}(AA^T) = \mathbf{rank}(A)$$

We now go forward when $A^T A$ is invertible to solve the normal equations $A^T A\hat{x} = A^T b$.

The Normal Equations $A^T A\hat{x} = A^T b$

Figure II.2 shows a picture of the least squares problem and its solution. The problem is that b is not in the column space of A , so $Ax = b$ has no solution. The best vector $p = A\hat{x}$ is a projection. **We project b onto the column space of A .** The vectors \hat{x} and $p = A\hat{x}$ come from solving a famous system of linear equations : $A^T A\hat{x} = A^T b$. To invert $A^T A$, we need to know that A has independent columns.

The picture shows the all-important right triangle with sides b, p , and e .

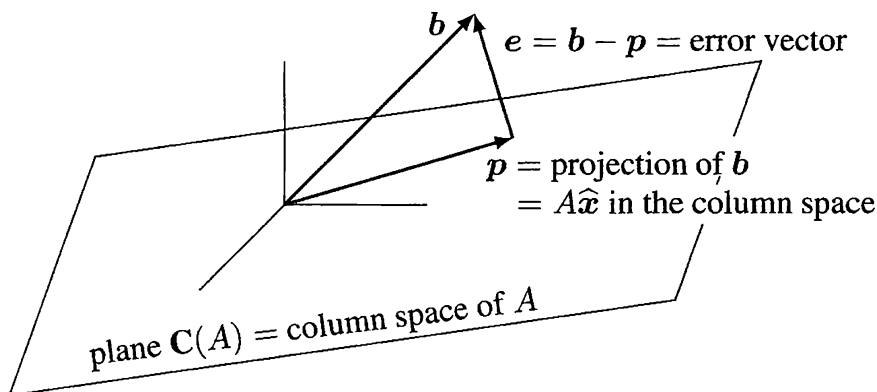


Figure II.2: The projection $p = A\hat{x}$ is the point in the column space that is closest to b .

Everybody understands that e is perpendicular to the plane (the column space of A). This says that $b - p = b - A\hat{x}$ is perpendicular to all vectors Ax in the column space :

$$(Ax)^T (b - A\hat{x}) = x^T A^T (b - A\hat{x}) = 0 \text{ for all } x \text{ forces } A^T (b - A\hat{x}) = 0. \quad (4)$$

Everything comes from that last equation, when we write it as $A^T A\hat{x} = A^T b$.

Normal equation for \hat{x}	$A^T A\hat{x} = A^T b$	(5)
Least squares solution to $Ax = b$	$\hat{x} = (A^T A)^{-1} A^T b$	(6)
Projection of b onto the column space of A	$p = A\hat{x} = A(A^T A)^{-1} A^T b$	(7)
Projection matrix that multiplies b to give p	$P = A(A^T A)^{-1} A^T$	(8)

A now has independent columns: $r = n$. That makes $A^T A$ positive definite and invertible. We could check that our \hat{x} is the same vector $x^+ = A^+ b$ that came from the pseudoinverse. There are no other \hat{x} 's because the rank is assumed to be $r = n$. The nullspace of A only contains the zero vector.

Projection matrices have the special property $P^2 = P$. When we project a second time, the projection p stays exactly the same. Use equation (8) for P :

$$P^2 = A(A^T A)^{-1} A^T A(A^T A)^{-1} A^T = A(A^T A)^{-1} A^T = P. \tag{9}$$

The Third Way to Compute \hat{x} : Gram-Schmidt

The columns of A are still assumed to be independent: $r = n$. But they are not assumed to be orthogonal! Then $A^T A$ is not a diagonal matrix and solving $A^T A \hat{x} = A^T b$ needs work. Our third approach **orthogonalizes the columns of A** , and then \hat{x} is easy to find.

You could say: The work is now in producing orthogonal (even orthonormal) columns. Exactly true. The operation count is actually doubled compared to $A^T A \hat{x} = A^T b$, but orthogonal vectors provide numerical stability. Stability becomes important when $A^T A$ is nearly singular. The **condition number** of $A^T A$ is its norm $\|A^T A\|$ times $\|(A^T A)^{-1}\|$. When this number σ_1^2/σ_n^2 is large, it is wise to orthogonalize the columns of A in advance. Then work with an orthogonal matrix Q .

The condition number of Q is $\|Q\|$ times $\|Q^{-1}\|$. *Those norms equal 1: best possible.*

Here is the famous Gram-Schmidt idea, starting with A and ending with Q . Independent columns a_1, \dots, a_n lead to orthonormal q_1, \dots, q_n . This is a fundamental computation in linear algebra. The first step is $q_1 = a_1/\|a_1\|$. That is a unit vector: $\|q_1\| = 1$. Then subtract from a_2 its component in the q_1 direction:

Gram-Schmidt step	Orthogonalize	$A_2 = a_2 - (a_2^T q_1) q_1$	(10)
	Normalize	$q_2 = A_2/\ A_2\ $	(11)

Subtracting that component $(a_2^T q_1) q_1$ produced the vector A_2 orthogonal to q_1 :

$$(a_2 - (a_2^T q_1) q_1)^T q_1 = a_2^T q_1 - a_2^T q_1 = 0 \text{ since } q_1^T q_1 = 1.$$

The algorithm goes onward to a_3 and A_3 and q_3 , normalizing each time to make $\|q\| = 1$. Subtracting the components of a_3 along q_1 and q_2 leaves A_3 :

Orthogonalize $A_3 = a_3 - (a_3^T q_1) q_1 - (a_3^T q_2) q_2$ **Normalize** $q_3 = \frac{A_3}{\|A_3\|}$ (12)

$$A_3^T q_1 = A_3^T q_2 = 0 \text{ and } \|q_3\| = 1$$

Each q_k is a combination of a_1 to a_k . Then each a_k is a combination of q_1 to q_k .

a_1	$= \ a_1\ q_1$	
a_2	$= (a_2^T q_1) q_1 + \ A_2\ q_2$	(13)
a_3	$= (a_3^T q_1) q_1 + (a_3^T q_2) q_2 + \ A_3\ q_3$	

Those equations tell us that **the matrix $R = Q^T A$ with $r_{ij} = q_i^T a_j$ is upper triangular :**

$$\begin{bmatrix} a_1 & a_2 & a_3 \end{bmatrix} = \begin{bmatrix} q_1 & q_2 & q_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{bmatrix} \text{ is } A = QR. \quad (14)$$

Gram-Schmidt produces orthonormal q 's from independent a 's. Then $A = QR$.

If $A = QR$ then $R = Q^T A =$ inner products of q 's with a 's! Later a 's are not involved in earlier q 's, so R is triangular. And certainly $A^T A = R^T Q^T Q R = R^T R$:

The least squares solution to $Ax = b$ is $\hat{x} = R^{-1} Q^T b$.

The MATLAB command is $[Q, R] = qr(A)$. Every $r_{ij} = q_i^T a_j$ because $R = Q^T A$. The vector $\hat{x} = (A^T A)^{-1} A^T b$ is $(R^T R)^{-1} R^T Q^T b$. This is exactly $\hat{x} = R^{-1} Q^T b$.

Gram-Schmidt with Column Pivoting

That straightforward description of Gram-Schmidt worked with the columns of A in their original order a_1, a_2, a_3, \dots . This could be dangerous! We could never live with a code for elimination that didn't allow row exchanges. Then roundoff error could wipe us out.

Similarly, each step of Gram-Schmidt should begin with a new column that is as independent as possible of the columns already processed. **We need column exchanges to pick the largest remaining column. Change the order of columns as we go.**

To choose correctly from the remaining columns of A , we make a simple change in Gram-Schmidt:

Old Accept column a_j as next. Subtract its components in the directions q_1 to q_{j-1}

New When q_{j-1} is found, subtract the q_{j-1} component from **all remaining columns**

This might look like more work, but it's not. Sooner or later we had to remove $(a_i^T q_{j-1})q_{j-1}$ from each remaining column a_i . Now we do it sooner—as soon as we know q_{j-1} . Then we have a free choice of the next column to work with, and we choose the largest.

Elimination Row exchanges on A left us with $PA = LU$ (permutation matrix P)

Gram-Schmidt Column exchanges leave us with $AP = QR$ (permutation matrix P)

*Here is the situation after $j - 1$ Gram-Schmidt steps with column pivoting. We have $j - 1$ orthogonal unit vectors q_1 to q_{j-1} in the columns of a matrix Q_{j-1} . We have the square matrix R_{j-1} that combines those columns of Q_{j-1} to produce $j - 1$ columns of A . They might not be the first $j - 1$ columns of A —we are optimizing the column order. All the remaining columns of A have been *orthogonalized* against the vectors q_1 to q_{j-1} .*

Step j . Choose the largest of the remaining columns of A . Normalize it to length 1.

This is q_j . Then from each of the $n - j$ vectors still waiting to be chosen, subtract the component in the direction of this latest q_j . Ready now for step $j + 1$.

We will follow Gunnar Martinsson's 2016 course notes for APPM 5720, to express step j in pseudocode. The original A is A_0 and the matrices Q_0 and R_0 are empty.

Step j is the following loop, which starts with A_{j-1} and ends at A_j . The code stops after j reaches $\min(m, n)$. **Here is column pivoting in Gram-Schmidt:**

$i = \operatorname{argmax} \|A_{j-1}(:, \ell)\|$ finds the largest column not yet chosen for the basis
 $\mathbf{q}_j = A_{j-1}(:, i) / \|A_{j-1}(:, i)\|$ normalizes that column to give the new unit vector \mathbf{q}_j
 $Q_j = [Q_{j-1} \quad \mathbf{q}_j]$ updates Q_{j-1} with the new orthogonal unit vector \mathbf{q}_j
 $\mathbf{r}_j = \mathbf{q}_j^T A_{j-1}$ finds the row of inner products of \mathbf{q}_j with remaining columns of A
 $R_j = \begin{bmatrix} R_{j-1} \\ \mathbf{r}_j \end{bmatrix}$ updates R_{j-1} with the new row of inner products
 $A_j = A_{j-1} - \mathbf{q}_j \mathbf{r}_j$ subtracts the new rank-one piece from each column to give A_j

When this loop ends, we have Q and R and $A = QR$. This R is a *permutation* of an upper triangular matrix. (It will be upper triangular if the largest columns in Step 1 come first, so each $i = j$.) The actual output can be an upper triangular matrix plus a vector with the numbers $1, \dots, n$ in the permutation order we need to know, to construct R .

In practice, this QR algorithm with pivoting is made safer by *reorthonormalizing*:

$$\mathbf{q}_j = \mathbf{q}_j - Q_{j-1}(Q_{j-1}^T \mathbf{q}_j)$$

$$\mathbf{q}_j = \mathbf{q}_j / \|\mathbf{q}_j\| \quad (\text{to make sure!})$$

There is a similar reordering for “ QR with Householder matrices” to reduce roundoff error. You have seen the essential point of pivoting: **good columns come first.**

Question: Both Q from Gram-Schmidt and U from the SVD contain an orthonormal basis for the column space $\mathbf{C}(A)$. *Are they likely to be the same basis?*

Answer: No, they are not the same. The columns of U are eigenvectors of AA^T . You cannot find eigenvectors (or eigenvalues) in a finite number of exact “arithmetic” steps for matrices of size $n > 4$. The equation $\det(A - \lambda I) = 0$ will be 5th degree or higher: *No formula can exist for the roots λ of a 5th degree equation (Abel).* Gram-Schmidt just requires inner products and square roots so Q must be different from U .

In the past, computing a nearly accurate eigenvalue took a much larger multiple of n^3 floating-point operations than elimination or Gram-Schmidt. That is not true now.