

Intro to Least Squares Problems (LSQ)

Suppose we want to find a solution to a system $\mathbf{Ax} = \mathbf{b}$, (with \mathbf{A} $m \times n$) that doesn't have an exact solution – maybe \mathbf{A} isn't full rank, or $m > n$, or for whatever reason \mathbf{b} isn't in the column-space of \mathbf{A} , i.e. $\mathbf{b} \notin \text{col}(\mathbf{A})$.

Then an alternative kind of “solution” is a vector \mathbf{x} that minimizes the 2-norm of the residual error, i.e. $\text{argmin}_{\mathbf{x}} \|\mathbf{b} - \mathbf{Ax}\|_2$. Here is a summary of approaches to finding such a solution. There are different cases, depending on the dimensions and rank of \mathbf{A} .

Full rank, $m = n$

Then \mathbf{A} is invertible, and so $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ gives a unique solution: $\mathbf{A}^{-1}\mathbf{b}$ is unique and $\|\mathbf{b} - \mathbf{Ax}\|_2 = 0$ is minimal since norms must be nonnegative.

Full rank, $m > n$

We explore a few approaches to this case. To motivate learning more about QR decompositions, we first explored using a QR factorization to solve the LSQ problem in this case.

Using QR decomposition

Let $\mathbf{A} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}} = [\mathbf{Q} \ \mathbf{U}] \begin{bmatrix} \mathbf{R} \\ 0 \end{bmatrix} = \mathbf{QR}$ for \mathbf{Q} an $m \times n$ matrix of orthonormal columns and \mathbf{R} an $n \times n$ upper-triangular matrix. Since \mathbf{A} is full rank, \mathbf{R} is invertible. Now let's use the orthogonal invariance of the 2-norm!

$$\begin{aligned} \|\mathbf{b} - \mathbf{Ax}\|_2 &= \|\tilde{\mathbf{Q}}^T (\mathbf{b} - \tilde{\mathbf{Q}}\tilde{\mathbf{R}}\mathbf{x})\|_2 \quad \text{orthogonal invariance of 2-norm, plugging in QR} \\ &= \left\| \begin{bmatrix} \mathbf{Q}^T \mathbf{b} \\ \mathbf{U}^T \mathbf{b} \end{bmatrix} - \begin{bmatrix} \mathbf{R}\mathbf{x} \\ 0 \end{bmatrix} \right\|_2 \\ &= \left\| \begin{bmatrix} \mathbf{Q}^T \mathbf{b} - \mathbf{R}\mathbf{x} \\ \mathbf{U}^T \mathbf{b} \end{bmatrix} \right\|_2 \end{aligned}$$

which leads to choosing $\mathbf{x} = \mathbf{R}^{-1}$ so cancel out the upper partition in the expression, yielding $\|\mathbf{U}^T \mathbf{b}\|_2$ as the minimized error.

Uniqueness of solution In class we touched on the question “is QR unique, or could we improve this error by choosing a different QR factorization that would yield smaller $\|\mathbf{U}^T \mathbf{b}\|_2$?” This is where I made the mistake of thinking \mathbf{U}^T was a tall narrow set of orthogonal columns, and I *incorrectly* stated $\|\mathbf{U}^T \mathbf{b}\|_2 = \|\mathbf{b}\|_2$ via the orthogonal invariance of the 2-norm.

Despite my mistake, the error *is* uniquely minimized: looking at the uniqueness of the QR decomposition (described in the [QR uniqueness](#) Lecture Appendix) we have that any other QR decomposition would yield a factor \mathbf{U}_2 related to our original \mathbf{U} by some square orthogonal matrix \mathbf{T} (which depends on \mathbf{U}_2) as follows: $\mathbf{U}_2 = \mathbf{UT}$. This allows us to write, for any choice of QR decomposition we might make, $\|\mathbf{U}_2^T \mathbf{b}\|_2 = \|(\mathbf{UT})^T \mathbf{b}\|_2 = \|\mathbf{T}^T \mathbf{U}^T \mathbf{b}\|_2 = \|\mathbf{U}^T \mathbf{b}\|_2$ using

the orthogonal invariance of the 2-norm – which this time we *can* use since T is *square* orthogonal.

This proves that the solution $\mathbf{x} = \mathbf{R}^{-1}\mathbf{Q}^T\mathbf{b}$ is the unique 2-norm minimizer.

Using the Normal Equations

We want to find \mathbf{x} that minimizes $\|\mathbf{b} - \mathbf{Ax}\|_2$; the same \mathbf{x} minimizes

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{b} - \mathbf{Ax}\|_2^2 \\ &= (\mathbf{b} - \mathbf{Ax})^T(\mathbf{b} - \mathbf{Ax}) \\ &= \mathbf{b}^T\mathbf{b} - \mathbf{x}^T\mathbf{A}^T\mathbf{b} - \mathbf{b}^T\mathbf{Ax} + \mathbf{x}^T\mathbf{A}^T\mathbf{Ax} \\ &= \mathbf{b}^T\mathbf{b} - 2\mathbf{b}^T\mathbf{Ax} + \mathbf{x}^T\mathbf{A}^T\mathbf{Ax}. \end{aligned}$$

The \mathbf{x} that minimizes $f(\mathbf{x})$ is the same \mathbf{x} that minimizes $h(\mathbf{x}) = -2\mathbf{b}^T\mathbf{Ax} + \mathbf{x}^T\mathbf{A}^T\mathbf{Ax}$ because $\mathbf{b}^T\mathbf{b}$ is a constant independent of \mathbf{x} . The minimum of the convex function $h(\mathbf{x})$ is where its gradient, $\nabla(\mathbf{x})$, is 0. we have

$$\frac{\partial}{\partial \mathbf{x}}(-2\mathbf{b}^T\mathbf{Ax}) = -2\mathbf{A}^T\mathbf{b} \quad \text{and} \quad \frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T\mathbf{A}^T\mathbf{Ax}) = 2\mathbf{A}^T\mathbf{Ax}$$

so $\nabla(\mathbf{x}) = -2\mathbf{A}^T\mathbf{b} + 2\mathbf{A}^T\mathbf{Ax}$. Setting $\nabla(\mathbf{x}) = 0$ yields the **normal equations**:

$$\mathbf{A}^T\mathbf{Ax} = \mathbf{A}^T\mathbf{b}. \tag{1}$$

Since we are in the case that \mathbf{A} is full rank and $m \times n$ with $m > n$, we have that $\mathbf{A}^T\mathbf{A}$ is invertible. Hence, $\mathbf{x} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$ is the unique minimizer of $h(\mathbf{x})$, thus the unique minimizer of $f(\mathbf{x})$, and hence of $\|\mathbf{b} - \mathbf{Ax}\|_2$.

(Side note: to see more on *why* a zero of the gradient of $h(\mathbf{x})$ is a minimizer of the function $h(\mathbf{x})$, take CS 520 on optimization. The intuition is the same reason that minima (and maxima) of a function of *one* variable, $f(x)$, are where the derivative is 0, i.e. $f'(x) = 0$.)

Using the SVD

Let $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ be the SVD for \mathbf{A} , where \mathbf{U} and \mathbf{V} are square, orthogonal. Then using the orthogonal invariance of the 2-norm we have

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{b} - \mathbf{Ax}\|_2^2 \\ &= \|\mathbf{U}^T(\mathbf{b} - \mathbf{U}\Sigma\mathbf{V}^T(\mathbf{V}\mathbf{V}^T)\mathbf{x})\|_2^2 \\ &= \|(\mathbf{U}^T\mathbf{b}) - \Sigma(\mathbf{V}^T\mathbf{x})\|_2^2. \end{aligned}$$

Since \mathbf{A} is $m \times n$ with $m > n$, we know $\Sigma = \begin{bmatrix} \mathbf{D} \\ 0 \end{bmatrix}$ where \mathbf{D} is $n \times n$ diagonal, $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n)$. From this, we can write

$$\begin{aligned} f(\mathbf{x}) &= \|(\mathbf{U}^T\mathbf{b}) - \Sigma(\mathbf{V}^T\mathbf{x})\|_2^2 \\ &= \left\| \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \end{bmatrix} - \begin{bmatrix} \mathbf{D} \\ 0 \end{bmatrix} (\mathbf{V}^T\mathbf{x}) \right\|_2^2 \end{aligned}$$

and so minimizing $f(\mathbf{x})$ is achieved by \mathbf{x} satisfying $\mathbf{V}^T\mathbf{x} = \mathbf{D}^{-1}\mathbf{c}_1$ where $\begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \end{bmatrix} = \mathbf{U}^T\mathbf{b}$. Multiplying by \mathbf{V} yields $\mathbf{x} = \mathbf{V}\mathbf{D}^{-1}\mathbf{c}_1 = \sum_{i=1}^n \frac{1}{\sigma_i} \mathbf{v}_i(\mathbf{u}_i^T\mathbf{b})$.

Since \mathbf{A} is full rank, we have $\sigma_i \neq 0$ for $i = 1 : n$, and so this all makes sense. Further note that, since the SVD is unique (up to factors of ± 1 on the vectors \mathbf{u}_i and \mathbf{v}_i), this expression of \mathbf{x} is unique.

In fact, using the expression below for \mathbf{A}^+ we have that

$$\begin{aligned} \mathbf{x} &= \mathbf{V}\mathbf{D}^{-1}\mathbf{c}_1 \\ &= \sum_{i=1}^n \frac{1}{\sigma_i} \mathbf{v}_i (\mathbf{u}_i^T \mathbf{b}) \\ &= \mathbf{A}^+ \mathbf{b} \\ &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} \end{aligned}$$

Generalized Inverses and the Rank-deficient LSQ problem

In the SVD approach we took $\Sigma = \begin{bmatrix} \mathbf{D} \\ 0 \end{bmatrix}$ where $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n)$ is an $n \times n$ diagonal matrix, and since we were in the case that \mathbf{A} had full rank, $\sigma_n \neq 0$, so \mathbf{D} was invertible. In class we discussed using this to define a “pseudo-inverse” for \mathbf{A} as follows:

$$\mathbf{A}^+ := \mathbf{V}\Sigma^+ \mathbf{U}^T = \mathbf{V} \begin{bmatrix} \mathbf{D}^{-1} & 0 \end{bmatrix} \mathbf{U}^T$$

where $\mathbf{D}^{-1} = \text{diag}(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_n})$ exists because we had assumed \mathbf{A} was full rank.

Now we abandon the assumption that \mathbf{A} is full rank. To solve this case we expand our definition of \mathbf{A}^+ to the case that $\sigma_i = 0$ for $i > r$, where $r = \text{rank}(\mathbf{A}) < n$. In this situation, $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$, and we define $\mathbf{D}^+ = \text{diag}(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r}, 0, \dots, 0)$. Then we can define the **generalized inverse** of $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ to be

$$\mathbf{A}^+ := \mathbf{V}\Sigma^+ \mathbf{U}^T = \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1} & & 0 & \dots & 0 \\ & \ddots & & & \\ & & \frac{1}{\sigma_r} & & 0 \\ & & & & 0 \\ & & & & 0 \end{bmatrix} \mathbf{U}^T = \sum_{j=1}^r \frac{1}{\sigma_j} \mathbf{v}_j \mathbf{u}_j^T.$$

This is one approach to the generalized inverse. Another approach is to define the generalized inverse to be the unique matrix satisfying certain properties:

Moore Penrose generalized inverse of an arbitrary $m \times n$ matrix \mathbf{A} is the unique $n \times m$ matrix \mathbf{X} satisfying

1. $\mathbf{A}\mathbf{X}\mathbf{A} = \mathbf{A}$
2. $\mathbf{X}\mathbf{A}\mathbf{X} = \mathbf{X}$
3. $(\mathbf{A}\mathbf{X})^T = \mathbf{A}\mathbf{X}$
4. $(\mathbf{X}\mathbf{A})^T = \mathbf{X}\mathbf{A}$

Remarks If \mathbf{A}^{-1} exists, then $\mathbf{X} = \mathbf{A}^{-1}$. The matrix \mathbf{A}^+ as defined above satisfies these properties, and so is equivalent to the Moore-Penrose generalized inverse. If \mathbf{A} is full rank and $m \geq n$ then we have $\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$. Finally, $\mathbf{A}\mathbf{A}^+$ and $\mathbf{A}^+ \mathbf{A}$ are the orthogonal projections of \mathbf{A} onto $\text{range}(\mathbf{A})$ and $\text{range}(\mathbf{A}^T)$, respectively.

Rank deficient, $m \geq n$

The rank deficient case follows the use of the SVD in the full rank case: with the generalized inverse defined in for even rank-deficient \mathbf{A} , the minimum norm solution to the LSQ problem is $\mathbf{x} = \mathbf{A}^+\mathbf{b}$.

Numerical Rank Since we solve these problems on computers, the exact rank of a matrix is not always the same as the rank of the floating-point representation of a matrix. For example, A matrix $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ with σ_j very small might have $\sigma_j = 0$ in floating-point arithmetic. To make this notion precise, we define the **numerical rank** of a matrix \mathbf{A} by $r(\varepsilon) = \#$ of singular values of \mathbf{A} that are $\geq \varepsilon$.