

Least-Squares Approximation and Bilinear Forms

The Normal Equations :

Suppose a column vector g is given in an Euclidean space into which a given matrix F maps another real space of column vectors u . In the Euclidean space, length $\|g\| := \sqrt{(g^T g)}$. Usually matrix F is rectangular with more rows than columns. Our task is to choose u to minimize $\|Fu - g\|$, which will then be the distance from g to $\text{Range}(F)$. This task is called “Least-Squares” because the u we seek would minimize the sum of squares of the elements of $Fu - g$.

Differentiating the sum of squares produces $d((Fu - g)^T(Fu - g)) = 2(Fu - g)^T F du$, which vanishes for all (infinitesimal) perturbations du if and only if $(Fu - g)^T F = 0^T$. (Do you see why?) Transposed, this becomes the “Normal Equations” of the Least-Squares Problem,

$$(F^T F)u = F^T g;$$

u must satisfy them to minimize $\|Fu - g\|$ and so best approximate g by a vector in $\text{Range}(F)$.

Do the Normal Equations have at least one solution $u = \hat{u}$?

If so, is $\|F\hat{u} - g\| \leq \|Fu - g\|$ for all u ? *i.e.*, does \hat{u} minimize rather than maximize?

These questions among others are addressed in this note.

At first sight one might think the Normal Equations’ solution should be $\hat{u} = (F^T F)^{-1} F^T g$. But this formula fails if the columns of F are linearly dependent. To see why, observe that

$$“Fz = 0” \iff “(Fz)^T(Fz) = z^T(F^T F)z = 0” \iff “(F^T F)z = 0”,$$

so $F^T F$ is invertible (nonsingular) if and only if the columns of F are linearly independent. This hasn’t been assumed; in fact matrix F could be rectangular with more columns than rows.

Exercise: Show that, if the rows of F are linearly independent, a solution $\hat{u} = F^T(F F^T)^{-1} g$ and, if not the only solution of the Normal Equations, it is the only solution that minimizes $\hat{u}^T \hat{u}$ too.

In short, if neither the rows nor the columns of F are linearly independent, neither $F F^T$ nor $F^T F$ need be invertible, and then the existence of a minimizing solution $u = \hat{u}$ is in question.

Warning: Even when an indicated inverse exists, neither formula $\hat{u} = (F^T F)^{-1} F^T g$ nor $\hat{u} = F^T(F F^T)^{-1} g$ should be used with numerical data unless the computer’s arithmetic carries at least twice as many sig. digits as are trusted in the data $[F, g]$ or desired in the result \hat{u} . Otherwise roundoff will degrade the result \hat{u} too badly whenever F is too near a matrix of lower rank. The reason behind this warning will become clear after *Singular Values* have been discussed. If the arithmetic carries barely more sig. digits than are trusted in the data or desired in the result, it should be computed by means of a *QR factorization*, which will also be discussed later. Matlab uses such a factorization to compute \hat{u} , which Matlab calls “ $F \backslash g$ ”, whenever F is not square. Least-Squares is built into Matlab.

Existence and Uniqueness of a Minimizing Solution \hat{u} :

We shall use Fredholm’s Alternatives (*q.v.*) to deduce that the Normal Equations always have at least one solution \hat{u} , and to determine when it is unique. At least one solution exists if and only if $w^T(F^T g) = 0$ whenever $w^T(F^T F) = 0^T$, so consider any row w^T that satisfies the last equation. It must satisfy also $0 = w^T(F^T F)w = (Fw)^T(Fw)$, which implies $Fw = 0$, which implies $w^T(F^T g) = (Fw)^T g = 0$, whereupon Fredholm’s Alternative (1) implies that the Normal Equations have at least one solution \hat{u} . It is unique if and only if the columns of F are linearly independent; otherwise add any nonzero solution z of $Fz = 0$ to one \hat{u} to get another.

How do we know that setting $u = \hat{u}$ minimizes $\|Fu - g\|$? For every u we find

$$\begin{aligned} \|Fu - g\|^2 - \|F\hat{u} - g\|^2 &= \|F(u - \hat{u}) + (F\hat{u} - g)\|^2 - \|F\hat{u} - g\|^2 \\ &= \|F(u - \hat{u})\|^2 + 2(F(u - \hat{u}))^T(F\hat{u} - g) \quad (\text{since } \|z\|^2 = z^T z) \\ &= \|F(u - \hat{u})\|^2 + 2(u - \hat{u})^T F^T(F\hat{u} - g) = \|F(u - \hat{u})\|^2 \geq 0, \end{aligned}$$

with equality instead of inequality just when u is a(nother) solution of the Normal Equations.

When the Normal Equations have many solutions \hat{u} , which does Matlab choose for $F \setminus g$? It has a near minimal number of nonzero elements. A different solution minimizes $\|\hat{u}\|^2 := \hat{u}^T \hat{u}$, as if also the space of vectors u were Euclidean. This doubly minimizing solution \hat{u} satisfies both the Normal Equations $(F^T F)\hat{u} = F^T g$ and an auxiliary equation $\hat{u} = F^T F v$ for some vector v of “Lagrange Multipliers.” In consequence v satisfies $(F^T F)^2 v = F^T g$, an equation with least one solution v whose existence is assured by an application of Fredholm’s first alternative very much like before except that the hypothesis $w^T(F^T F) = 0^T$ is replaced by $w^T(F^T F)^2 = 0^T$. (Can you carry out this inference?) Every other solution u of the Normal Equations satisfies $\|u\|^2 - \|\hat{u}\|^2 = \|(u - \hat{u}) + \hat{u}\|^2 - \|\hat{u}\|^2 = \dots = \|u - \hat{u}\|^2 + 2\hat{u}^T(u - \hat{u}) = \|u - \hat{u}\|^2 + 2v^T F^T F(u - \hat{u}) = \|u - \hat{u}\|^2 \geq 0$, so this $\hat{u} = F^T F v$ really is doubly minimizing; moreover it is determined uniquely by the data $[F, g]$. (Can you see why?) As we shall see later after Singular Values have been discussed, there is a matrix F^\dagger called the “Moore-Penrose Pseudo-Inverse” of F such that the doubly minimizing $\hat{u} = F^\dagger g$ is a linear function of g . (Matlab’s name for F^\dagger is `pinv(F)`.) However, whenever neither $F^T F$ nor FF^T is invertible, so F^\dagger is interesting, it turns out to be a violently discontinuous function of F . This renders the doubly minimizing \hat{u} doubly dubious because the space of vectors u need not be Euclidean. Matlab’s $F \setminus g$ can be discontinuous too, even when FF^T is invertible and the doubly minimizing \hat{u} is continuous.

Linear Regression:

Least-Squares approximation has been applied to statistical estimation for over two centuries. An m -by- n matrix F is assumed given with linearly independent columns (so $m \geq n$); and a given m -vector $g = y + q$ of “data” is thought to include a systematic contribution y and a “random error” q . The question is how near is y to $\text{Range}(F)$? The answer is obscured by the random error. The elements of this error q are assumed *independently distributed* with *mean* 0 and known *variance* β^2 . These terms are given meaning by an *Averaging* or *Expectation* operator \mathcal{A} which acts upon every random variable r linearly to produce $\mathcal{A}r$, the average or *mean* of the population of values of r . Thus $\mathcal{A}q = 0$ because every element of q has mean 0; and q has *covariance* matrix $\mathcal{A}((q - \mathcal{A}q)(q - \mathcal{A}q)^T) = \beta^2 I$ since the square of every element of q has mean β^2 but every product of different elements of q has mean 0 because they are independent. The smaller is β , the less uncertainty does random error q introduce into the data g .

Define $x := (F^T F)^{-1} F^T y$ to minimize $\|Fx - y\|$ although neither y nor x can be known. As the known g approximates y , so is x approximated by whatever \hat{u} minimizes $\|F\hat{u} - g\|$. Get

$\hat{u} = (F^T F)^{-1} F^T g$; how well can it approximate x ? Since $\mathcal{A}g = y$, we find that $\mathcal{A}\hat{u} = x$, so \hat{u} is an *unbiased* estimate of x . The covariance matrix of \hat{u} is computable too; it is

$\mathcal{A}((\hat{u} - x)(\hat{u} - x)^T) = \mathcal{A}((F^T F)^{-1} F^T q q^T F (F^T F)^{-1}) = (F^T F)^{-1} F^T \mathcal{A}(q q^T) F (F^T F)^{-1} = \beta^2 (F^T F)^{-1}$. The smaller this is, the better does \hat{u} approximate x on average. The smaller is $\|Fx - y\|$, the smaller do we expect $\|F\hat{u} - g\|$ to be. How small should we expect it to be? A calculation below shows that $\mathcal{A}(\|F\hat{u} - g\|^2) = \|Fx - y\|^2 + (m-n)\beta^2$. It means that $\|F\hat{u} - g\|$ is unlikely to exceed $\beta\sqrt{(m-n)}$ much if y lies in or very near $\text{Range}(F)$; conversely, $\|Fx - y\|$ is unlikely to be much smaller than $\|F\hat{u} - g\|$ if this is many times bigger than $\beta\sqrt{(m-n)}$. Explanation follows.

Proof that $\mathcal{E}(\|F\hat{u} - g\|^2) = \|Fx - y\|^2 + (m-n)\beta^2$: The Trace of a square matrix is defined to be the sum of its diagonal elements; evaluate this sum to confirm that $\text{Trace}(B^T C) = \text{Trace}(C B^T)$ for any matrices B^T and C whose products $B^T C$ and $C B^T$ are both square, though perhaps of different dimensions. Next define $H := F(F^T F)^{-1} F^T$ and confirm that $H^T = H = H^2$. (H is the orthogonal projector onto $\text{Range}(F)$ because “ $p = Fz$ for some z ” \iff “ $p = Hp$ ”, so $\text{Range}(F) = \text{Range}(H)$, and “ $H z = o$ ” \iff “ $z^T H = o^T$ ”, so $\text{Nullspace}(H) = \text{Range}(H)^\perp$.) Shortly we shall have use for $\text{Trace}(H) = \text{Trace}((F^T F)^{-1} F^T F) = \text{Trace}(I_n) = n$. Now we observe that \hat{u} and x are so defined that $F\hat{u} - g = (H - I)g$ and $Fx - y = (H - I)y$ wherein I is the m -by- m identity matrix. Consequently

$$\begin{aligned} \mathcal{E}(\|F\hat{u} - g\|^2) &= \mathcal{E}(((H-I)g)^T (H-I)g) = \mathcal{E}(\text{Trace}((H-I)g((H-I)g)^T)) \quad \dots \text{ because } \text{Trace}(b^T c) = \text{Trace}(c b^T) \\ &= \mathcal{E}(\text{Trace}((H-I)g g^T (H-I))) = \text{Trace}((H-I)\mathcal{E}(g g^T)(H-I)) \quad \dots \text{ because } H = H^T \text{ isn't random} \\ &= \text{Trace}((H-I)\mathcal{E}(y y^T + y q^T + q y^T + q q^T)(H-I)) \quad \dots \text{ because } g = y + q \\ &= \text{Trace}((H-I)(y y^T + O + O + \beta^2 I)(H-I)) = \text{Trace}((H-I)y y^T (H-I)) + \beta^2 \text{Trace}((H-I)^2) \\ &= \text{Trace}((Fx - y)(Fx - y)^T) + \beta^2 \text{Trace}(I - H) = \|Fx - y\|^2 + \beta^2(m-n) \quad \text{as was claimed.} \end{aligned}$$

Proof that $\|F\hat{u} - g\|^2$ is unlikely to be many times bigger than its mean $\mathcal{E}(\|F\hat{u} - g\|^2)$: More precisely, we shall deduce that $\|F\hat{u} - g\|^2$ exceeds $\lambda \mathcal{E}(\|F\hat{u} - g\|^2)$ with probability less than $1/\lambda$ for every $\lambda > 1$. This deduction is an instance of *Tchebyshev's Inequality*: If a positive random variable ρ has mean $\mu := \mathcal{E}\rho$, then the probability that $\rho \geq \lambda\mu$ cannot exceed $1/\lambda$ for any $\lambda > 1$. Here is a proof of Tchebyshev's Inequality. Let $p(\xi)$ be the probability that $\rho \leq \xi$. This $p(\xi)$ is a nondecreasing function increasing from $p(0) = 0$ to $p(\infty) = 1$, and $\mu = \int_0^\infty \xi \, d p(\xi)$ by virtue of the definition of \mathcal{E} . We seek an overestimate for $\int_{\lambda\mu}^\infty d p(\xi)$, which is the probability that $\rho \geq \lambda\mu$. We find that $\int_{\lambda\mu}^\infty d p(\xi) \leq \int_{\lambda\mu}^\infty \xi \, d p(\xi) / (\lambda\mu) \leq \int_0^\infty \xi \, d p(\xi) / (\lambda\mu) = \mu / (\lambda\mu)$, which yields the result claimed. (This can be a gross overestimate because it uses almost no information about p . For almost all values of $\lambda > 1$, and for all values of $\lambda > 1$ for almost all probability functions p , the probability that $\rho \geq \lambda\mu$ is actually far tinier than $1/\lambda$.)

Thus the computed $\|F\hat{u} - g\|^2$ is unlikely to be many times bigger than $\|Fx - y\|^2 + \beta^2(m-n)$ in which $\beta^2(m-n)$ is given and $\|Fx - y\|^2$ is unknown, whence something probabilistic can be inferred about the unknown. Another similar application of Least-Squares is to the assumption that $y = Fx$ and $g = y + q$ for a random error q about which β^2 is unknown but estimated from $\|F\hat{u} - g\|^2 / (m-n)$. These applications are treated in Statistics courses.

Abstract Least-Squares:

Suppose a column vector g is given in an Euclidean space into which a given linear operator \mathbf{F} maps a real space of abstract vectors \mathbf{u} . In the Euclidean space, length $\|g\| := \sqrt{(g^T g)}$, but no such length is defined (yet) for $\text{Domain}(\mathbf{F})$. Again our task is to choose \mathbf{u} to minimize $\|\mathbf{F}\mathbf{u} - g\|$, which will then be the distance from g to $\text{Range}(\mathbf{F})$. Differentiating the sum of squares $\|\mathbf{F}\mathbf{u} - g\|^2 = (\mathbf{F}\mathbf{u} - g)^T (\mathbf{F}\mathbf{u} - g)$ produces $d((\mathbf{F}\mathbf{u} - g)^T (\mathbf{F}\mathbf{u} - g)) = 2(\mathbf{F}\mathbf{u} - g)^T \mathbf{F} \, d\mathbf{u}$, which vanishes for all (infinitesimal) perturbations $d\mathbf{u}$ if and only if $(\mathbf{F}\mathbf{u} - g)^T \mathbf{F} = \mathbf{o}^T$. This \mathbf{o}^T is the linear functional that annihilates $\text{Domain}(\mathbf{F})$. The last equation says that when $\|\mathbf{F}\mathbf{u} - g\|$ is minimized the residual $\mathbf{F}\mathbf{u} - g$ must be *normal* (perpendicular, orthogonal) to $\text{Range}(\mathbf{F})$. (This explains the word “Normal” in “Normal Equations” and removes any suggestion that other equations are abnormal.) Drawing a picture helps; imagine $\text{Range}(\mathbf{F})$ to be a plane in Euclidean 3-space containing a vector $\mathbf{F}\mathbf{u}$ which, when it comes closest to a given vector g not in the plane, comes to that point in the plane reached by dropping a perpendicular from g .

We could transpose “ $(\mathbf{F}\mathbf{u} - g)^T \mathbf{F} = \mathbf{o}^T$ ” to “ $(\mathbf{F}^T \mathbf{F})\mathbf{u} = \mathbf{F}^T g$ ” if we knew what “ $\mathbf{F}^T \mathbf{F}$ ” meant.

The trouble with the expression “ $\mathbf{F}^T\mathbf{F}$ ” is that it is not what it first seems; if \mathbf{F} were a matrix then $\mathbf{F}^T\mathbf{F}$ would map $\text{Domain}(\mathbf{F})$ to itself, but a change of basis in $\text{Domain}(\mathbf{F})$ does not change $\mathbf{F}^T\mathbf{F}$ to the expected *similar* matrix. Here is what happens instead:

Let \mathbf{B} be a basis for $\text{Domain}(\mathbf{F})$. Then abstract vector $\mathbf{u} = \mathbf{B}\mathbf{u}$ for some column vector \mathbf{u} , and $\mathbf{F}\mathbf{u} = \mathbf{F}\mathbf{B}\mathbf{u} = \mathbf{F}\mathbf{u}$ for a matrix $\mathbf{F} = \mathbf{F}\mathbf{B}$. The Normal Equations “ $(\mathbf{F}\mathbf{u} - \mathbf{g})^T\mathbf{F} = \mathbf{o}^T$ ” turn into “ $(\mathbf{F}\mathbf{u} - \mathbf{g})^T\mathbf{F} = \mathbf{o}^T$ ” which becomes “ $(\mathbf{F}^T\mathbf{F})\mathbf{u} = \mathbf{F}^T\mathbf{g}$ ” after matrix transposition. $\mathbf{B}\mathbf{C}$ is a new basis for $\text{Domain}(\mathbf{F})$, and $\mathbf{u} = \mathbf{B}\mathbf{C}\mathbf{u}$ for $\mathbf{u} = \mathbf{C}^{-1}\mathbf{u}$, and $\mathbf{F}\mathbf{u} = \mathbf{F}\mathbf{u}$ for matrix $\mathbf{F} = \mathbf{F}\mathbf{C}$, where \mathbf{C} is any invertible matrix of the same dimension as $\text{Domain}(\mathbf{F})$. What was “ $(\mathbf{F}^T\mathbf{F})\mathbf{u} = \mathbf{F}^T\mathbf{g}$ ” in the old basis becomes “ $(\mathbf{F}^T\mathbf{F})\mathbf{u} = \mathbf{F}^T\mathbf{g}$ ” in the new, replacing matrix $\mathbf{F}^T\mathbf{F}$ by $\mathbf{F}^T\mathbf{F} = \mathbf{C}^T\mathbf{F}^T\mathbf{F}\mathbf{C}$. This differs from $\mathbf{C}^{-1}\mathbf{F}^T\mathbf{F}\mathbf{C}$, which is how the change in basis would have changed $\mathbf{F}^T\mathbf{F}$ if it were the matrix of a map from $\text{Domain}(\mathbf{F})$ to itself. Instead, $\mathbf{F}^T\mathbf{F}$ is the matrix of a map from $\text{Domain}(\mathbf{F})$ to its own dual space.

If you doubt that these choices of basis matter, try the following example: Let $\mathbf{g} := 10101$, a scalar, and suppose $\mathbf{F} = [1, 10, 100]$ in some coordinate system. Then get Matlab to compute $\mathbf{u} = \mathbf{F} \setminus \mathbf{g}$ to solve the least-squares problem. Next change to a new basis using a diagonal matrix $\mathbf{C} = \text{diag}([10, 1, 1/16])$. It changes \mathbf{F} to $\mathbf{F} = \mathbf{F}\mathbf{C}$ and thus changes the solution of the least-squares problem to $\mathbf{u} = \mathbf{F} \setminus \mathbf{g}$. This maps back to $\mathbf{C}\mathbf{u} = \mathbf{C} * ((\mathbf{F} * \mathbf{C}) \setminus \mathbf{g})$ in the old basis. Compare with the old solution \mathbf{u} . Try again with 6-vectors \mathbf{g} and 6-by-3 matrices \mathbf{F} at random.

Bilinear Forms:

There is no uniquely defined operator $\mathbf{F}^T\mathbf{F}$ just as there is no functional \mathbf{u}^T determined uniquely by vector \mathbf{u} in a non-Euclidean space. The matrices that appear in the Normal Equations are not all matrices that represent linear maps from one space of column vectors to another or itself; matrix $\mathbf{F}^T\mathbf{F}$ belongs to a *Symmetric Bilinear Form* that maps column vectors to row vectors.

Consider $(\mathbf{F}\mathbf{u})^T\mathbf{F}\mathbf{v}$. It maps pairs $\{\mathbf{u}, \mathbf{v}\}$ of vectors from $\text{Domain}(\mathbf{F})$ to real scalars, and does so as a linear function of each vector separately; this is the definition of a *Bilinear Form*. And since $(\mathbf{F}\mathbf{u})^T\mathbf{F}\mathbf{v}$ is unaltered when \mathbf{u} and \mathbf{v} are swapped, it is a *Symmetric Bilinear Form*.

There are many notations for bilinear forms: $\mathbf{H}\mathbf{u}\mathbf{v}$, $H(\mathbf{u}, \mathbf{v})$, $(\mathbf{v}, \mathbf{H}\mathbf{u})$, They all mean this:

$\mathbf{H}\mathbf{u}$ is a linear functional in the space dual to vectors \mathbf{v} , and $\mathbf{H}\mathbf{u}\mathbf{v}$ is its scalar value;

$\mathbf{H}\mathbf{v}$ is a linear functional in the space dual to vectors \mathbf{u} , and $\mathbf{H}\mathbf{u}\mathbf{v}$ is its scalar value;

Given a basis \mathbf{B} for vectors $\mathbf{u} = \mathbf{B}\mathbf{u}$, and a basis \mathbf{E} for vectors $\mathbf{v} = \mathbf{E}\mathbf{v}$, there is a

matrix \mathbf{H} for which $\mathbf{H}\mathbf{u}\mathbf{v} = (\mathbf{H}\mathbf{B}\mathbf{u})\mathbf{E}\mathbf{v} = (\mathbf{H}\mathbf{u})^T\mathbf{v} = \mathbf{v}^T\mathbf{H}\mathbf{u}$;

Changing bases from \mathbf{B} to $\mathbf{B}\mathbf{C}$ and \mathbf{E} to $\mathbf{E}\mathbf{D}$ changes \mathbf{u} to $\mathbf{u} = \mathbf{C}^{-1}\mathbf{u}$, \mathbf{v} to $\mathbf{v} = \mathbf{D}^{-1}\mathbf{v}$,

and \mathbf{H} to $\mathbf{H} = \mathbf{D}^T\mathbf{H}\mathbf{C}$ so that $\mathbf{H}\mathbf{u}\mathbf{v} = \mathbf{v}^T\mathbf{H}\mathbf{u} = \mathbf{v}^T\mathbf{H}\mathbf{u}$.

Exercise: Express the elements of matrix \mathbf{H} in terms of the effect \mathbf{H} has upon the elements of bases \mathbf{B} and \mathbf{E} .

A *Symmetric* bilinear form maps vectors \mathbf{u} and \mathbf{v} from the same space to scalars, and does so in a way independent of the order of \mathbf{u} and \mathbf{v} thus: $\mathbf{H}\mathbf{u}\mathbf{v} = \mathbf{H}\mathbf{v}\mathbf{u}$. A symmetric bilinear form has a symmetric matrix $\mathbf{H} = \mathbf{H}^T$ in any basis. (Why?) Changing the basis changes \mathbf{H} to matrix $\mathbf{H} = \mathbf{C}^T\mathbf{H}\mathbf{C}$ for some invertible \mathbf{C} ; the two matrices \mathbf{H} and \mathbf{H} are called “Congruent.” This congruence is an *Equivalence*, so it preserves rank; *i.e.*, $\text{rank}(\mathbf{H}) = \text{rank}(\mathbf{H})$. Congruence also preserves a thing called “Signature” as we’ll see when we come to Sylvester’s *Inertia* Theorem.