

LEAST SQUARE PROBLEM

Consider observation data $b_1, b_2, \dots, b_m \in \mathbb{R}$ of some quantities, e.g., temperature, test scores, Dow Jones index, accidents, you name it. Meanwhile, for each observation b_i , we pair with some independent data $a_{i1}, a_{i2}, \dots, a_{in}$, such as, humidity, age, unemployment rate, eyesight, etc.. An example: we observe that at 32 Fahrenheit (b_i , dependent variable), local humidity is 5%, local wind speed is 32 mph, and local pressure is 1 atm. In fact, we may draw up data in a table.

Temperature in F \mathbf{b}	Humidity \mathbf{a}_1	Wind Speed \mathbf{a}_2	Pressure \mathbf{a}_3	Air Quality Index \mathbf{a}_4
32	5	17	1	55
37	7	13	1.02	65
44	8	11	0.99	14
47	11	6	0.96	36

We seek the answer to the question: how does temperature depend on all these variables with the provided data? It is completely natural to posit a linear relationship between \mathbf{b} and the \mathbf{a}_i 's. More precisely, for each b_i , we seek the coefficients x_1, \dots, x_n that

$$b_i = x_0 + x_1 a_{i1} + x_2 a_{i2} + \dots + x_n a_{in}, \quad i = 1, 2, \dots, m.$$

Is it utterly possible that we can find the exact x_{i1}, \dots, x_{in} that satisfy this relationship? We probably need n equations at least to determine these unknowns. So,

$$\begin{aligned} b_1 &= x_0 + x_1 a_{11} + \dots + x_n a_{1n} = (1, a_{11}, \dots, a_{1n}) \cdot (x_0, x_1, \dots, x_n) \\ &\dots \\ &\dots \\ b_m &= x_0 + x_1 a_{m1} + \dots + x_n a_{mn} = (1, a_{m1}, \dots, a_{mn}) \cdot (x_0, x_1, \dots, x_n) \end{aligned}$$

or more compactly,

$$\mathbf{b} = \begin{bmatrix} 1 & a_{11} & a_{12} & \dots & a_{1n} \\ 1 & a_{21} & \dots & \dots & \dots \\ 1 & \dots & \dots & \dots & \dots \\ 1 & \dots & \dots & \dots & \dots \\ 1 & \dots & \dots & \dots & \dots \\ 1 & a_{m1} & \dots & \dots & a_{mn} \end{bmatrix}_{m \times (n+1)} \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ \dots \\ x_n \end{bmatrix}_{(n+1) \times 1} = \mathbf{Ax}.$$

So, we have $(n + 1)$ unknowns but only m equations. This is only possibly exactly solvable when $m = n + 1$ which means we have exactly the same number of observation data points as the number of independent variables. This is unrealistic. In practice, $m \gg n + 1$, that is, we have massive amount of data, but only a handful of features we seek the extent of dependence on. Therefore, the system $\mathbf{Ax} = \mathbf{b}$ here is not solvable in general.

Then what? Game over? We go for the next best thing. Now, if $\mathbf{Ax} = \mathbf{b}$ is not solvable, we may find some \mathbf{x} that achieves $\mathbf{Ax} \approx \mathbf{b}$, which is a reasonable request. In fact, we seek a solution that minimizes the square of the l^2 -error

$$\phi(\mathbf{x}) = \|\mathbf{b} - \mathbf{Ax}\|_2^2$$

where $\phi: \mathbb{R}^{n+1} \rightarrow \mathbb{R}$. Recall that $\mathbf{b} - \mathbf{Ax}$ is known as the residual vector. All we are trying to do is to come up with a solution that minimizes the l^2 -norm of this residual vector.

We expand the l^2 -norm by definition,

$$\phi(\mathbf{x}) = \phi(x_0, x_1, \dots, x_n) = \sum_{i=1}^m (b_i - x_0 - x_1 a_{i1} - x_2 a_{i2} - \dots - x_n a_{in})^2.$$

How do we minimize a function of multiple variables? We compute its gradient and set it equal to $\mathbf{0}$ to find the critical points.

$$0 = \frac{\partial \phi}{\partial x_j} = -2 \sum_{i=1}^m (b_i - x_0 - x_1 a_{i1} - x_2 a_{i2} - \dots - x_n a_{in}) a_{ij}, \quad j = 0, 1, \dots, n.$$

Moving the -2 out of the way, we see that the critical point(s) satisfy

$$\sum_{i=1}^m (b_i - x_0 - x_1 a_{i1} - x_2 a_{i2} - \cdots - x_n a_{in}) a_{ij} = 0, \quad j = 0, 1, \dots, n.$$

Guess what? Now, we have exactly $n + 1$ equations for the $n + 1$ unknowns. This set of equations is called the **normal equations**.

Denote $y_i = b_i - x_0 - x_1 a_{i1} - x_2 a_{i2} - \cdots - x_n a_{in}$ and $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$. Then, the equation reads

$$\sum_{i=1}^m a_{ij} y_i = 0, \quad j = 0, 1, \dots, n.$$

which now requires you to recall the definition of matrix vector multiplication –

$$(\mathbf{Ax})_i = (a_{i1}, a_{i2}, \dots, a_{in}) \cdot (x_1, x_2, \dots, x_n) \implies i^{\text{th}} \text{ row of } \mathbf{A} \text{ dotted with } \mathbf{x}.$$

Let's visualize what $\sum_{i=1}^m a_{ij} y_i$ really is:

$$(a_{1j}, a_{2j}, a_{3j}, \dots, a_{mj}) \cdot (y_1, \dots, y_m)$$

where we realize that $(a_{1j}, a_{2j}, a_{3j}, \dots, a_{mj})$ is the j^{th} **column** of \mathbf{A} , which means it is the j^{th} **row** of \mathbf{A}^T . Thus,

$$\sum_{i=1}^m a_{ij} y_i = (\mathbf{A}^T \mathbf{y})_j.$$

Enumerating over all $j = 1, 2, \dots, m$, we find that the **normal equations** can be written in matrix form,

$$\mathbf{A}^T \mathbf{y} = \mathbf{0}.$$

Now, looking at the definition of \mathbf{y} , we have

$$\begin{aligned} y_1 &= b_1 - x_0 - x_1 a_{11} - x_2 a_{12} - \cdots - x_n a_{1n} \\ &\cdot \\ &\cdot \\ y_m &= b_m - x_0 - x_1 a_{m1} - x_2 a_{m2} - \cdots - x_n a_{mn} \end{aligned}$$

which is

$$\mathbf{y} = \mathbf{b} - \mathbf{Ax}.$$

Inserting this back into the normal equation $\mathbf{A}^T \mathbf{y} = \mathbf{0}$, we have

$$\mathbf{A}^T (\mathbf{b} - \mathbf{Ax}) = \mathbf{0} \implies \mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b},$$

the celebrated final form of the **normal equations**. All we need is the observation data: the dependent variable \mathbf{b} , and the independent variables \mathbf{A} .

We put the problem in full form: the (minimizer) solution to the **least square problem**

$$\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{b} - \mathbf{Ax}\|_2^2$$

is the solution to the set of **normal equations**

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

Now, after finding where the critical point is, we still need to confirm that this critical point indeed gives me the minimum, not the maximum.

Theorem. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Every solution $\tilde{\mathbf{x}}$ to $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$ satisfies

$$\|\mathbf{b} - \mathbf{Ax}\|_2 \leq \|\mathbf{b} - \mathbf{Ax}\|_2 \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

that is, $\tilde{\mathbf{x}}$, if exists, is the global minimizer of $\|\mathbf{b} - \mathbf{Ax}\|_2$.

Proof. Given $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$, we have

$$\|\mathbf{u} + \mathbf{v}\|_2^2 = (\mathbf{u} + \mathbf{v})^T (\mathbf{u} + \mathbf{v}) = \|\mathbf{u}\|_2^2 + 2\mathbf{u}^T \mathbf{v} + \|\mathbf{v}\|_2^2.$$

Then,

$$\begin{aligned} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 &= \|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}} + \mathbf{A}\tilde{\mathbf{x}} - \mathbf{A}\mathbf{x}\|_2^2 \\ &= \|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2 + 2(\mathbf{A}(\tilde{\mathbf{x}} - \mathbf{x}))^T (\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}) + \|\mathbf{A}(\tilde{\mathbf{x}} - \mathbf{x})\|_2^2 \\ &\geq \|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2 + 2(\tilde{\mathbf{x}} - \mathbf{x})^T \mathbf{A}^T (\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}) \xrightarrow{0} \\ &= \|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2. \end{aligned}$$

□

EXISTENCE OF A SOLUTION

It remains to show that $\tilde{\mathbf{x}}$ indeed exists, and under one more condition on \mathbf{A} , is also unique. Existence is not hard if we know a little bit of linear algebra. Note that $\mathbf{A}^T \mathbf{b}$ lies in the range of \mathbf{A}^T . But we also can show that the range of \mathbf{A}^T and that of $\mathbf{A}^T \mathbf{A}$ are the same (a fundamental theorem in linear algebra), which means there exists \mathbf{x} such that $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b}$ since both sides of the equation maps to the same subspace.

UNIQUENESS OF THE SOLUTION

If $\det(\mathbf{A}^T \mathbf{A}) \neq 0$, then we are all set because then $\mathbf{A}^T \mathbf{A}$ is invertible, and

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

But is this always the case? This should depend on \mathbf{A} – but here \mathbf{A} is not necessarily square. So the usual technique from linear algebra won't work.

The claim here is that \mathbf{A} must have linearly independent columns iff $\mathbf{A}^T \mathbf{A}$ is invertible. For the forward direction, assume that \mathbf{A} has linearly independent columns, we suppose, on the contrary, that $\mathbf{A}^T \mathbf{A}$ is not invertible. Then, $\det(\mathbf{A}^T \mathbf{A}) = 0$, which means there exists nonzero $\mathbf{z} = \mathbf{0}$ such that

$$\mathbf{A}^T \mathbf{A}\mathbf{z} = \mathbf{0}.$$

Now, multiplying \mathbf{z}^T on the left, we have

$$\mathbf{z}^T \mathbf{A}^T \mathbf{A}\mathbf{z} = 0 \implies (\mathbf{A}\mathbf{z})^T (\mathbf{A}\mathbf{z}) = 0 \implies \|\mathbf{A}\mathbf{z}\|_2^2 = 0 \implies \mathbf{A}\mathbf{z} = \mathbf{0}, \quad \text{where } \mathbf{z} \neq \mathbf{0}.$$

But this is impossible because \mathbf{A} has linearly independent columns, i.e., the only solution to $\mathbf{A}\mathbf{z} = \mathbf{0}$ is $\mathbf{z} = \mathbf{0}$. Contradiction!