

Solving the least squares problem

We now return to the general least squares optimization problem¹ of

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2, \quad (1)$$

where \mathbf{A} is an $M \times N$ matrix and $\mathbf{y} \in \mathbb{R}^M$.

In the context of linear regression, where $N = 2$, we were able to derive a 2×2 system of equations that we could then solve to obtain a formula for the solution to the least squares problem. Here we would like to do something similar for the general case.

Recall that our approach was based on computing the partial derivative of the objective function with respect to the two parameters we were trying to estimate. If we want to take the same approach in the general setting, we have N parameters to consider, and thus, there are N partial derivatives to set to zero. In this case, the natural way to organize our computations is using the **gradient**. Given a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$, we denote the gradient (with respect to \mathbf{x}) of f by

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_N} \end{bmatrix}.$$

The gradient is just the vector in \mathbb{R}^N of partial derivatives with all N components of \mathbf{x} . If your multivariable calculus is rusty, see the brief refresher at the end of these notes.

¹Note that while this is the same problem as last time, we have exchanged \mathbf{x} in the place of $\boldsymbol{\alpha}$.

We are now tempted to solve the least squares problem by computing the gradient of the objective function in (1) and setting it equal to zero to obtain a system of equations that will yield the solution to the least squares problem. In particular, one can show (this is worked out as an example at the end of these notes) that this gradient is given by

$$\nabla \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{y})$$

Setting this equal to zero (or more specifically, a vector of zeros) yields the system

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

If $\mathbf{A}^T \mathbf{A}$ is invertible, then simply inverting this gives us the formula

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}.$$

You may recognize the formula $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$, which is one possible formula for what is called the **pseudo-inverse** of \mathbf{A} , and is often denoted by \mathbf{A}^\dagger . Provided that $\mathbf{A}^T \mathbf{A}$ is in fact invertible, this does indeed provide the optimal least squares estimate, i.e., the solution to (1).

We will discuss the properties of this estimate in much more detail later, including a more careful discussion about what happens if $\mathbf{A}^T \mathbf{A}$ is not actually invertible, but first I want to provide a bit more justification as to why taking the gradient and setting it equal to zero does in fact give us a solution to (1). Later in the course we will generalize this argument to arbitrary **convex** functions, but it is particularly simple in the least squares problem.

For convenience, let $f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$. Recall the fact that for any (column) vector $\mathbf{z} \in \mathbb{R}^N$, we can write $\|\mathbf{z}\|_2^2 = \mathbf{z}^T \mathbf{z}$. Thus, we can

also write

$$\begin{aligned}
 f(\mathbf{x}) &= (\mathbf{y} - \mathbf{A}\mathbf{x})^\top (\mathbf{y} - \mathbf{A}\mathbf{x}) \\
 &= \mathbf{y}^\top \mathbf{y} - (\mathbf{A}\mathbf{x})^\top \mathbf{y} - \mathbf{y}^\top \mathbf{A}\mathbf{x} + (\mathbf{A}\mathbf{x})^\top \mathbf{A}\mathbf{x} \\
 &= \mathbf{y}^\top \mathbf{y} - \mathbf{x}^\top \mathbf{A}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{A}\mathbf{x} + \mathbf{x}^\top \mathbf{A}^\top \mathbf{A}\mathbf{x} \\
 &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{y} + \mathbf{x}^\top \mathbf{A}^\top \mathbf{A}\mathbf{x},
 \end{aligned}$$

where above we use the following facts:

- For any matrices \mathbf{A} and \mathbf{B} , $(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$.
- For any matrices \mathbf{A} and \mathbf{B} , $(\mathbf{A}\mathbf{B})^\top = \mathbf{B}^\top \mathbf{A}^\top$.
- Using the previous fact, we have $\mathbf{x}^\top \mathbf{A}^\top \mathbf{y} = (\mathbf{y}^\top \mathbf{A}\mathbf{x})^\top$. Since these are scalars, we also have $\mathbf{y}^\top \mathbf{A}\mathbf{x} = (\mathbf{y}^\top \mathbf{A}\mathbf{x})^\top$, and hence $\mathbf{x}^\top \mathbf{A}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{A}\mathbf{x}$.

With this in hand, we consider a small perturbation $f(\mathbf{x} + \mathbf{u})$. Our goal will be to show that if $\hat{\mathbf{x}}$ corresponds to a vector where $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$, then \mathbf{x} must be a solution to (1), which is equivalent to showing that

$$f(\hat{\mathbf{x}}) \leq f(\hat{\mathbf{x}} + \mathbf{u})$$

for any possible choice of \mathbf{u} .

To do this, note that following the same argument as above we can write

$$\begin{aligned}
 f(\hat{\mathbf{x}} + \mathbf{u}) &= (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}} - \mathbf{A}\mathbf{u})^\top (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}} - \mathbf{A}\mathbf{u}) \\
 &= (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}})^\top (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}) - 2(\mathbf{A}\mathbf{u})^\top (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}) + (\mathbf{A}\mathbf{u})^\top (\mathbf{A}\mathbf{u}) \\
 &= \|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_2^2 - \mathbf{u}^\top (2\mathbf{A}^\top (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}})) + \|\mathbf{A}\mathbf{u}\|_2^2
 \end{aligned}$$

Thus, we obtain

$$\begin{aligned} f(\hat{\mathbf{x}} + \mathbf{u}) &= f(\hat{\mathbf{x}}) + \mathbf{u}^T(2\mathbf{A}^T(\mathbf{A}\hat{\mathbf{x}} - \mathbf{y})) + \|\mathbf{A}\mathbf{u}\|_2^2 \\ &\geq f(\hat{\mathbf{x}}) + \mathbf{u}^T(2\mathbf{A}^T(\mathbf{A}\hat{\mathbf{x}} - \mathbf{y})), \end{aligned}$$

where the last inequality follows since (by the definition of a norm) we have $\|\mathbf{A}\mathbf{u}\|_2^2 \geq 0$. Now, recall that earlier we showed that $\nabla f(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{y})$. Thus, if $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$, then from the above we have

$$f(\hat{\mathbf{x}} + \mathbf{u}) \geq f(\hat{\mathbf{x}}) + \mathbf{u}^T \nabla f(\hat{\mathbf{x}}) = f(\hat{\mathbf{x}}).$$

Since this holds for any possible choice of \mathbf{u} , this establishes that $\hat{\mathbf{x}}$ is indeed the minimizer of f and the solution to (1).

To summarize:

If the matrix $\mathbf{A}^T \mathbf{A}$ is invertible, then the optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$$

has solution

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}.$$

Review of multivariable calculus

Recall that for a differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$, the **derivative** can be defined as

$$f'(x) = \lim_{\delta \rightarrow 0} \frac{f(x + \delta) - f(x)}{\delta}.$$

When dealing with a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ that is defined on N -dimensional vectors, we can define the **partial derivative** with respect to x_n as

$$\frac{\partial f(\mathbf{x})}{\partial x_n} = \lim_{\delta \rightarrow 0} \frac{f(\mathbf{x} + \delta \mathbf{e}_n) - f(\mathbf{x})}{\delta},$$

where \mathbf{e}_n is the n^{th} “standard basis element”, i.e., the vector of all zeros with a single 1 in the n^{th} entry.

The **gradient** of a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ can be viewed as the vector of partial derivatives given by:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_N} \end{bmatrix}.$$

We will use the term gradient in two subtly different ways. Sometimes we use $\nabla f(\mathbf{x})$ to describe a *vector-valued function* or a *vector field*, i.e., a function that takes an arbitrary $\mathbf{x} \in \mathbb{R}^N$ and produces another vector. However, we also use the term gradient, and the same notation $\nabla f(\mathbf{x})$, to refer to vector given by the gradient at a particular point \mathbf{x} . So sometimes when we say “gradient” we mean a

vector-valued function, and sometimes we mean a single vector, and in both cases we use the notation $\nabla f(\mathbf{x})$. Which one will usually be obvious by the context.²

Note that in some cases we will use the notation $\nabla_{\mathbf{x}} f(\mathbf{x})$ to indicate that we are taking the gradient with respect to \mathbf{x} . This can be helpful when f is a function of more variables than just \mathbf{x} , but most of the time this is not necessary so we will typically use the simpler $\nabla f(\mathbf{x})$.

The gradient is one of the most fundamental concepts of this course. We can interpret the gradient in many ways. One way to think of the gradient when evaluated at a particular point \mathbf{x} is that it defines a linear mapping from \mathbb{R}^N to \mathbb{R} . Specifically, given a $\mathbf{u} \in \mathbb{R}^N$, we can use $\nabla f(\mathbf{x})$ to define a mapping of \mathbf{u} to \mathbb{R} by simply taking the inner product between the two vectors:

$$\langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle.$$

What does this mapping tell us? One can show (although we will not prove this here) that it computes the **directional derivative** of f in the direction of \mathbf{u} , i.e.,

$$\langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle = \lim_{\delta \rightarrow 0} \frac{f(\mathbf{x} + \delta \mathbf{u}) - f(\mathbf{x})}{\delta}. \quad (2)$$

This tells us how fast f is changing at \mathbf{x} when we move in the direction of \mathbf{u} .

A related way to think of $\nabla f(\mathbf{x})$ is as a vector that is pointing in the direction of *steepest ascent*, i.e., the direction in which f increases the fastest when starting at \mathbf{x} . To justify this, note that we just

²This is just like in the scalar case, where the notation $f(x)$ can sometimes refer to the function f and sometimes the function evaluated at x .

observed that we can interpret $\langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle$ as measuring how quickly f increases when we move in the direction of \mathbf{u} . How can we find the direction \mathbf{u} that maximizes this quantity? You may recall that the Cauchy-Schwarz inequality tells us that

$$|\langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle| \leq \|\nabla f(\mathbf{x})\|_2 \|\mathbf{u}\|_2,$$

and that this holds with equality when \mathbf{u} is co-linear with $\nabla f(\mathbf{x})$, i.e., when \mathbf{u} points in the same direction as $\nabla f(\mathbf{x})$. Specifically, this implies that $\nabla f(\mathbf{x})$ is the direction of steepest *ascent*, and $-\nabla f(\mathbf{x})$ is the direction of steepest *descent*.

More broadly, this characterizes the entire sets of ascent/descent directions. Suppose that $f : \mathbb{R}^N \rightarrow \mathbb{R}$ is differentiable at \mathbf{x} . If $\mathbf{u} \in \mathbb{R}^N$ is a vector obeying $\langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle < 0$, then we say that \mathbf{u} is a **descent direction** from \mathbf{x} , and for small enough $t > 0$,

$$f(\mathbf{x} + t\mathbf{u}) < f(\mathbf{x}).$$

Similarly, if $\langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle > 0$, then we say that \mathbf{u} is an **ascent direction** from \mathbf{x} , and for small enough $t > 0$,

$$f(\mathbf{x} + t\mathbf{u}) > f(\mathbf{x}).$$

It should hopefully not be a huge stretch of the imagination to see that being able to compute the direction of steepest ascent (or steepest descent) will be useful in the context of finding a maximum/minimum of a function.

Examples:

1. Compute the gradient of $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$.

Note that

$$\frac{\partial f(\mathbf{x})}{\partial x_n} = \frac{\partial}{\partial x_n}(a_1x_1 + \cdots + a_Nx_N) = a_n,$$

and thus $\nabla f(\mathbf{x}) = \mathbf{a}$.

2. Compute the gradient of $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}$.

Note that

$$\frac{\partial f(\mathbf{x})}{\partial x_n} = \frac{\partial}{\partial x_n}(x_1^2 + \cdots + x_N^2) = 2x_n,$$

and thus $\nabla f(\mathbf{x}) = 2\mathbf{x}$.

One final note on the gradient. Here we adopt the convention that the gradient is a *column vector*. This is probably the most common choice and is most convenient in this class, but some texts will instead treat the gradient as a row vector. The reason for this is that this makes the gradient a special case of the **Jacobian**. For a vector-valued function $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$, the Jacobian is the matrix of partial derivatives

$$\mathbf{D}_f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_N} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_2(\mathbf{x})}{\partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial x_1} & \frac{\partial f_M(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_M(\mathbf{x})}{\partial x_N} \end{bmatrix}.$$

Note that if $f : \mathbb{R}^N \rightarrow \mathbb{R}$, then $\mathbf{D}_f(\mathbf{x}) = (\nabla f(\mathbf{x}))^\top$, i.e., the Jacobian is simply the gradient, but treated as a row vector.

We will also occasionally need to make use of second derivatives. For a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$, this is captured by the **Hessian**, which is

the matrix of all possible pairwise partial derivatives:

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_N} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_N \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_N \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_N \partial x_N} \end{bmatrix}.$$

Note that if we view the gradient $\nabla f(\mathbf{x})$ as a vector valued function mapping from \mathbb{R}^N to \mathbb{R}^N , then the Hessian is the same as the Jacobian of the gradient..

Example:

Compute the Hessian of $f(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$.

Recall that

$$\frac{\partial}{\partial x_n} f(\mathbf{x}) = 2x_n.$$

Thus

$$\frac{\partial^2}{\partial x_m \partial x_n} f(\mathbf{x}) = \begin{cases} 2 & \text{if } m = n, \\ 0 & \text{otherwise.} \end{cases}$$

Or, more compactly, $\nabla^2 f(\mathbf{x}) = 2\mathbf{I}$, where \mathbf{I} is the $N \times N$ identity matrix.

One final tool from multivariable calculus that will be incredibly useful in this course (and that, once you understand it, makes the algorithm for training a neural network seem incredibly obvious) is the **multivariate chain rule**. Specifically, imagine we have $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ and $g : \mathbb{R}^K \rightarrow \mathbb{R}^N$. Now consider the composition of these functions denoted by $f \circ g(\mathbf{x}) = f(g(\mathbf{x}))$. Suppose we want

to calculate the Jacobian (or gradient if $M = 1$) of $f \circ g$. The multivariate extension of the chain rule in this case tells us that:

$$\mathbf{D}_{f \circ g}(\mathbf{x}) = \mathbf{D}_f(g(\mathbf{x})) \cdot \mathbf{D}_g(\mathbf{x}).$$

Note that

- $f \circ g$ is a mapping from \mathbb{R}^K to \mathbb{R}^M , and $\mathbf{x} \in \mathbb{R}^K$.
- $\mathbf{D}_g(\mathbf{x})$ is an $N \times K$ matrix (the evaluation of the Jacobian of g at \mathbf{x}).
- $\mathbf{D}_f(g(\mathbf{x}))$ is an $M \times N$ matrix (the evaluation of the Jacobian of f at $g(\mathbf{x})$).
- $\mathbf{D}_{f \circ g}(\mathbf{x})$ is an $M \times K$ matrix, which matches what we expect since $f \circ g : \mathbb{R}^K \rightarrow \mathbb{R}^M$, and is also what we get by multiplying an $M \times N$ matrix with an $N \times K$ matrix.

Also note that if $M = 1$, this procedure returns a row vector. If we wish to report the gradient as a column vector, then we can do this simply by taking the transpose, so that

$$\nabla f \circ g = (\mathbf{D}_f(g(\mathbf{x})) \cdot \mathbf{D}_g(\mathbf{x}))^T = \mathbf{D}_g(\mathbf{x})^T \cdot \mathbf{D}_f(g(\mathbf{x}))^T.$$

Example:

Compute the gradient of $\|\mathbf{y} - \mathbf{Ax}\|_2^2$.

Define $h = \|\mathbf{y} - \mathbf{Ax}\|_2^2$. Note that $h = f \circ g$ where $f(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$ and $g(\mathbf{x}) = \mathbf{y} - \mathbf{Ax}$. We have already seen above that

$$\mathbf{D}_f(\mathbf{x}) = (\nabla f(\mathbf{x}))^T = 2\mathbf{x}^T.$$

It is also a simple consequence of the fact that $\nabla(\mathbf{a}^\top \mathbf{x}) = \mathbf{a}$ to show that

$$\mathbf{D}_g(\mathbf{x}) = -\mathbf{A}.$$

Putting this together gives us

$$\mathbf{D}_h(\mathbf{x}) = \mathbf{D}_f(g(\mathbf{x})) \cdot \mathbf{D}_g(\mathbf{x}) = 2(\mathbf{y} - \mathbf{Ax})^\top (-\mathbf{A}).$$

Using convention that gradient is a column vector, by taking the transpose and re-arranging, we have

$$\nabla \|\mathbf{y} - \mathbf{Ax}\|_2^2 = 2\mathbf{A}^\top (\mathbf{Ax} - \mathbf{y}).$$