# Some Vector Calculus Identities

Sohom Paul

February 9, 2021

**Fact 1.** $\nabla_x(a^\mathsf{T} x) = a$

*Proof.* Here, we have some intimidating notation (gradients and inner products!) hiding a really simple computation. Let's start with a concrete example. Suppose that $x = \langle u, v \rangle$ and $a = \langle 3, 4 \rangle$. Then $f(u, v) = 3u + 4v$, so $f_u = 3$ and $f_v = 4$ (where subscripts denote partial derivatives), so $\nabla f = \langle 3, 4 \rangle = a$. In general, we have

$$a^\mathsf{T} x = a_1 x_1 + a_2 x_2 + \ldots a_n x_n$$
$$\frac{\partial}{\partial x_i} a^\mathsf{T} x = \frac{\partial}{\partial x_i}\left(a_1 x_1 + a_2 x_2 + \ldots + a_i x_i + \ldots a_n x_n\right)$$
$$= 0 + 0 + \ldots + a_i + \ldots + 0$$
$$= a_i$$

Because the $i$th component of $\nabla_x(a^\mathsf{T} x)$ matches the $i$th component of $a$, the two vectors are equal, which is what we set out to prove. $\qquad\square$

**Fact 2.** $\nabla_x(x^\mathsf{T} A x) = (A + A^\mathsf{T})x$

*Proof.* Now, this compuation is a little bit trickier than the previous one. We notice that

$$x^\mathsf{T} A x = \begin{pmatrix} x_1 & x_2 & \ldots & x_n \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & \ldots & a_{1n} \\ a_{21} & a_{22} & \ldots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \ldots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$= \begin{pmatrix} x_1 & x_2 & \ldots & x_n \end{pmatrix} \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + \ldots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \ldots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \ldots + a_{nn}x_n \end{pmatrix}$$

$$= x_1 a_{11} x_1 + x_1 a_{12} x_2 \ldots + x_2 a_{21} x_1 + x_2 a_{22} x_2 + \ldots$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{n} x_i a_{ij} x_j$$

Now, to take the partial derivative with respect to $x_k$, there are 4 cases:

$$\frac{\partial}{\partial x_k}(x_i a_{ij} x_j) = \begin{cases} a_{kj} x_j & i = k, j \neq k \\ a_{ik} x_i & i \neq k, j = k \\ 2a_{kk} x_k & i = j = k \\ 0 & \text{otherwise} \end{cases}$$

It follows

$$\frac{\partial}{\partial x_k}(x^\mathsf{T} A x) = \frac{\partial}{\partial x_k}(x_i a_{ij} x_j)$$

$$= 2a_{kk} + \sum_{\substack{1 \leq j \leq n \\ j \neq k}} a_{kj} x_j + \sum_{\substack{1 \leq i \leq n \\ i \neq k}} a_{ik} x_i$$

$$= 2a_{kk} + \sum_{\substack{1 \leq j \leq n \\ j \neq k}} (a_{kj} + a_{jk}) x_j$$

$$= \sum_{1 \leq j \leq n} (a_{kj} + a_{jk}) x_j$$

$$= \begin{pmatrix} a_{k1} + a_{1k} & a_{k2} + a_{2k} & \ldots & a_{kn} + a_{nk} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$\nabla_x(x^\mathsf{T} A x) = \begin{pmatrix} a_{11} + a_{11} & a_{12} + a_{21} & \ldots & a_{1n} + a_{n1} \\ a_{21} + a_{12} & a_{22} + a_{22} & \ldots & a_{2n} + a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} + a_{1n} & a_{n2} + a_{2n} & \ldots & a_{nn} + a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$= (A + A^\mathsf{T}) x$$

which is what we set out to prove. $\qquad \square$

**Theorem 1.** $\arg\max_w \|Xw - y\|_2 = (X^\mathsf{T} X)^{-1}(X^\mathsf{T} y)$

*Proof.* You've likely seen a proof of this formula in your linear algebra class using some properties of orthogonal projections. I want to outline here an alternate proof using vector calculus that is easier to generalize to nonlinear models.

The key insight to recall the fact from elementary calculus that the derivative of a smooth function is 0 when it achieves a local minimum or maximum. The same holds for a function of multiple variables, except we set the gradient of the function to 0 (which is equivalent to setting all the partials to 0). Thus, minimizing the squared $L^2$ loss, we get

$$f(w) = \|Xw - y\|_2^2$$

$$= (Xw - y)^\mathsf{T}(Xw - y)$$

$$= (w^\mathsf{T} X^\mathsf{T} - y^\mathsf{T})(Xw - y)$$

$$= w^\mathsf{T} X^\mathsf{T} X w - y^\mathsf{T} X w - w^\mathsf{T} X^\mathsf{T} y + y^\mathsf{T} y$$

$$= w^\mathsf{T} X^\mathsf{T} X w - 2(X^\mathsf{T} y)^\mathsf{T} w + \|y\|_2^2$$

$$\nabla_w f(w) = (X^\mathsf{T} X + (X^\mathsf{T} X)^\mathsf{T}) w - 2X^\mathsf{T} y$$

$$= 2(X^\mathsf{T} X) w - 2X^\mathsf{T} y$$

where we used Facts 1 and 2 when taking the gradient. Setting the last line to 0 immediately gives the optimal coefficient vector as $w^* = (X^\mathsf{T}X)^{-1}X^\mathsf{T}y$. (It can be checked that $X^\mathsf{T}X$ is invertible if $X$ has linearly independent columns, which always holds in practice.) As this is the only point where the gradient is 0, we know that this is the global optimum. $\qquad\square$

Note that to a machine learning practitioner, our matrix $X$ has rows that correspond to *data points* and columns that correspond to *data features*. Our model is that the output is linear in the different features, with vector $w$ telling us how much weight we assign to each. The vector $y$ tells us the observed *labels*. Our goal is to fit some weight vector that minimizes the $L^2$ loss. As you will see in your homework, the "linear" in linear regression only refers to the fact that your model is a linear combination of features; by appropriately choosing data features, you can fit polynomials or other functions using linear regression. Neat!