

Notes on Andrew Ng's CS 229 Machine Learning Course

Tyler Neylon

331.2016

These are notes I'm taking as I review material from Andrew Ng's CS 229 course on machine learning. Specifically, I'm watching [these videos](#) and looking at the written notes and assignments posted [here](#). These notes are available in two formats: [html](#) and [pdf](#).

I'll organize these notes to correspond with the written notes from the class.

1 On lecture notes 1

The notes in this section are based on [lecture notes 1](#).

1.1 Gradient descent in general

Given a cost function $J(\theta)$, the general form of an update is

$$\theta_j := \theta_j - \alpha \frac{\partial J}{\partial \theta_j}.$$

It bothers me that α is an arbitrary parameter. What is the best way to choose this parameter? Intuitively, it could be chosen based on some estimate or actual value of the second derivative of J . What can be theoretically guaranteed about the rate of convergence under appropriate conditions?

Why not use Newton's method? A general guess: the second derivative of J becomes cumbersome to work with.

It seems worthwhile to keep my eye open for opportunities to apply improved optimization algorithms in specific cases.

1.2 Gradient descent on linear regression

I realize this is a toy problem because linear regression in practice is not solve iteratively, but it seems worth understanding well. The general update equation is, for a single example i ,

$$\theta_j := \theta_j + \alpha(y^{(i)} - h_\theta(x^{(i)}))x_j^{(i)}.$$

The delta makes sense in that it is proportional to the error $y - h_\theta$, and in that the sign of the product $(y - h_\theta)x$ guarantees moving in the right direction. However, my first guess would be that the expression $(y - h_\theta)/x$ would provide a better update.

For example, suppose we have a single data point (x, y) where $x \neq 0$, and a random value of θ . Then a great update would be

$$\theta_1 := \theta_0 + (y - \theta_0 x)/x,$$

since the next hypothesis value h_θ would then be

$$h_\theta = \theta_1 x = \theta_0 x + y - \theta_0 x = y,$$

which is good. Another intuitive perspective is that we should be making *bigger* changes to θ_j when x_j is *small*, since it's harder to influence h_θ for such x values.

This is not yet a solidified intuition. I'd be interested in revisiting this question if I have time.

1.3 Properties of the trace operator

The trace of a square matrix obeys the nice property that

$$\text{tr } AB = \text{tr } BA. \tag{1}$$

One way to see this is to note that

$$\text{tr } AB = a_{ij}b_{ji} = \text{tr } BA,$$

where I'm using the informal shorthand notation that a variable repeated within a single product implies that the sum is taken over all relevant values of that variable. Specifically,

$$a_{ij}b_{ji} \text{ means } \sum_{i,j} a_{ij}b_{ji}.$$

I wonder if there's a more elegant way to verify (1).

Ng gives other interesting trace-based equations, examined next.

- Goal: $\nabla_A \text{tr } AB = B^T$.

Since

$$\text{tr } AB = a_{ij}b_{ji},$$

we have that

$$(\nabla_A \text{tr } AB)_{ij} = b_{ji},$$

verifying the goal.

- Goal: $\nabla_{A^T} f(A) = (\nabla_A f(A))^T$.

This can be viewed as

$$(\nabla_{A^T} f(A))_{ij} = \frac{\partial f}{\partial a_{ji}} = (\nabla_A f(A))_{ji}.$$

- Goal: $\nabla_A \text{tr}(ABA^T C) = CAB + C^T AB^T$.

I'll use some nonstandard index variable names below because I think it will help avoid possible confusion. Start with

$$(ABA^T C)_{xy} = a_{xz}b_{zw}a_{vw}c_{vy}.$$

Take the trace of that to arrive at

$$\alpha = \text{tr}(ABA^T C) = a_{xz}b_{zw}a_{vw}c_{vx}.$$

Use the product rule to find $\frac{\partial \alpha}{\partial a_{ij}}$. You can think of this as, in the equation above, first setting $xz = ij$ for one part of the product rule output, and then setting $vw = ij$ for the other part. The result is

$$(\nabla_A \alpha)_{ij} = b_{jw}a_{vw}c_{vi} + a_{xz}b_{zj}c_{ix} = c_{vi}a_{vw}b_{jw} + c_{ix}a_{xz}b_{zj}.$$

(The second equality above is based on the fact that we're free to rearrange terms within products in the repeated-index notation being used. Such rearrangement is commutativity of numbers, not of matrices.)

This last expression is exactly the ij^{th} entry of the matrix $CAB + C^T AB^T$, which was the goal.