

Gradient Descent for Solving Least-Squares Problems

Proof: Bounds on Step Size for Guaranteed Convergence

©Barry Van Veen 2019

Gradient descent minimizes the cost function

$$f(\mathbf{w}) = \|\mathbf{Aw} - \mathbf{d}\|_2^2$$

using the iterative algorithm

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \tau \mathbf{A}^T(\mathbf{Aw}^{(k)} - \mathbf{d}), \quad k = 0, 1, 2, 3, \dots \quad (1)$$

where $\tau > 0$ so that we modify the current iterate in the negative gradient direction. Often this algorithm is initialized with $\mathbf{w}^{(0)} = \mathbf{0}$. The initialization does not affect the convergence behavior because $f(\mathbf{w})$ is convex.

The iteration is guaranteed to converge to the minimum of the cost function if the squared error decreases with each iteration, that is, if

$$f(\mathbf{w}^{(k+1)}) = \|\mathbf{Aw}^{(k+1)} - \mathbf{d}\|_2^2 < f(\mathbf{w}^{(k)}) = \|\mathbf{Aw}^{(k)} - \mathbf{d}\|_2^2$$

Begin by substituting $\mathbf{w}^{(k)} - \tau \mathbf{A}^T(\mathbf{Aw}^{(k)} - \mathbf{d})$ for $\mathbf{w}^{(k+1)}$ in $f(\mathbf{w}^{(k+1)})$ to write

$$f(\mathbf{w}^{(k+1)}) = \|\mathbf{A}(\mathbf{w}^{(k)} - \tau \mathbf{A}^T(\mathbf{Aw}^{(k)} - \mathbf{d})) - \mathbf{d}\|_2^2 \quad (2)$$

$$= \|(\mathbf{Aw}^{(k)} - \mathbf{d}) - \tau(\mathbf{AA}^T(\mathbf{Aw}^{(k)} - \mathbf{d}))\|_2^2 \quad (3)$$

Now let $\mathbf{c} = \mathbf{Aw}^{(k)} - \mathbf{d}$ and $\mathbf{e} = \tau(\mathbf{AA}^T(\mathbf{Aw}^{(k)} - \mathbf{d}))$ be the first and second terms in parentheses so $f(\mathbf{w}^{(k+1)}) = \|\mathbf{c} - \mathbf{e}\|_2^2 = (\mathbf{c} - \mathbf{e})^T(\mathbf{c} - \mathbf{e})$. Expand the product to write $f(\mathbf{w}^{(k+1)}) = \|\mathbf{c}\|_2^2 + \|\mathbf{e}\|_2^2 - 2\mathbf{e}^T\mathbf{c}$. Substituting for \mathbf{c} and \mathbf{e} we thus obtain

$$f(\mathbf{w}^{(k+1)}) = \|\mathbf{Aw}^{(k)} - \mathbf{d}\|_2^2 + \tau^2 \|\mathbf{AA}^T(\mathbf{Aw}^{(k)} - \mathbf{d})\|_2^2 - 2\tau ((\mathbf{Aw}^{(k)} - \mathbf{d})^T \mathbf{AA}^T)(\mathbf{Aw}^{(k)} - \mathbf{d}) \quad (4)$$

$$= f(\mathbf{w}^{(k)}) + \tau^2 \|\mathbf{A}(\mathbf{A}^T(\mathbf{Aw}^{(k)} - \mathbf{d}))\|_2^2 - 2\tau ((\mathbf{Aw}^{(k)} - \mathbf{d})^T \mathbf{A})(\mathbf{A}^T(\mathbf{Aw}^{(k)} - \mathbf{d})) \quad (5)$$

Define $\mathbf{v} = \mathbf{A}^T(\mathbf{Aw}^{(k)} - \mathbf{d})$ to simplify the expression and rewrite Eq. 5 as

$$f(\mathbf{w}^{(k+1)}) = f(\mathbf{w}^{(k)}) + \tau^2 \|\mathbf{Av}\|_2^2 - 2\tau \mathbf{v}^T \mathbf{v}$$

Note that \mathbf{v} does not depend on τ . Thus, to prove $f(\mathbf{w}^{(k+1)}) < f(\mathbf{w}^{(k)})$, we must find the condition for which

$$q(\tau) = \tau^2 \|\mathbf{Av}\|_2^2 - 2\tau \mathbf{v}^T \mathbf{v}$$

is less than zero.

Recall the operator norm of a matrix \mathbf{X} satisfies $\max_{\mathbf{g}} \|\mathbf{X}\mathbf{g}\|_2 \leq \|\mathbf{X}\|_{op} \|\mathbf{g}\|_2$, so the first term in $q(\tau)$ may be upper bounded as

$$\tau^2 \|\mathbf{A}\mathbf{v}\|_2^2 \leq \tau^2 \|\mathbf{A}\|_{op}^2 \|\mathbf{v}\|_2^2$$

We may rewrite the second term in $q(\tau)$ as

$$-2\tau \mathbf{v}^T \mathbf{v} = -2\tau \|\mathbf{v}\|_2^2$$

Hence, we obtain an upper bound on $q(\tau)$

$$q(\tau) \leq \tau^2 \|\mathbf{A}\|_{op}^2 \|\mathbf{v}\|_2^2 - 2\tau \|\mathbf{v}\|_2^2$$

Factoring out the common terms we write

$$q(\tau) \leq (\tau \|\mathbf{A}\|_{op}^2 - 2) \tau \|\mathbf{v}\|_2^2$$

The second term in $q(\tau)$, $\tau \|\mathbf{v}\|_2^2$, is positive provided $\mathbf{v} \neq \mathbf{0}$, so we obtain $q(\tau) < 0$ by requiring

$$(\tau \|\mathbf{A}\|_{op}^2 - 2) < 0$$

which indicates τ must satisfy

$$\tau < \frac{2}{\|\mathbf{A}\|_{op}^2}$$

Note that $\mathbf{v} = \mathbf{0}$ implies $\mathbf{A}^T(\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}) = \mathbf{0}$, or $\mathbf{A}^T \mathbf{A} \mathbf{w}^{(k)} = \mathbf{A}^T \mathbf{d}$, or $\mathbf{w}^{(k)} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{d}$. Thus, if $\mathbf{v} = \mathbf{0}$, then the iteration has converged to the minimum of the squared error and the update term in Eq. 1 is zero.

Hence, the gradient descent algorithm will converge to the minimum of the squared error cost function provided the step-size τ satisfies $\tau < \frac{2}{\|\mathbf{A}\|_{op}^2}$.