

ECEn 671: Mathematics of Signals and Systems

Randal W. Beard

Brigham Young University

September 1, 2023

Section 1

Gradient Descent

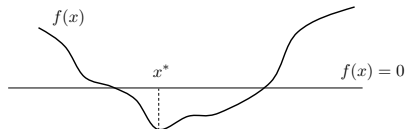
Gradient Descent

The topic for the remainder of the course is minimization and maximization of functions.

In particular we will constrain our attention to continuously differentiable functions.

Gradient Descent

Suppose we have a function of the form



and we would like to find x^* , what should we do?

Gradient Descent

The basic idea of gradient descent is to pick any $x^{[0]}$ and then move “downward”. To move down, we look at the slope of f .

If $\frac{\partial f}{\partial x}(x^{[k]})$ is positive, chose $x^{[k+1]} < x^{[k]}$.

If $\frac{\partial f}{\partial x}(x^{[k]})$ is negative, choose $x^{[k+1]} > x^{[k]}$

i.e.

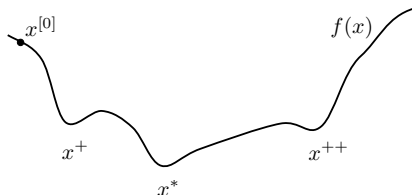
$$x^{[k+1]} = x^{[k]} - \alpha \frac{\partial f}{\partial x}(x^{[k]}),$$

where α is the step size.

Gradient Descent

Before moving to the multivariable case, let's consider the potential problems with this approach.

Problem 1: Local Minima. If f looks like this:



then if the initial condition is at $x^{[0]}$, the iteration

$$x^{[k+1]} = x^{[k]} - \alpha \frac{\partial f}{\partial x}(x^{[k]})$$

will converge to x^+ , if α is small enough.

Gradient Descent

Other initial conditions will result in x^{++} while others will give x^* , the true minimum.

This is a fundamental problem with any method that relies on derivative information. There are no completely satisfactory solutions to the problem. However there are many ad-hoc fixes.

Example

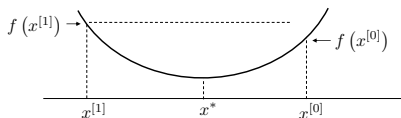
- ▶ Execute from numerous “random” initial conditions and pick the lowest solution.
- ▶ Occasionally introduce random jumps in x .
- ▶ etc...

Gradient Descent

Problem 2: Step Size. The selection of α can have a major effect on the convergence of the sequence

$$x^{[k+1]} = x^{[k]} - \alpha \frac{\partial f}{\partial x}(x^{[k]})$$

For example,



Note f is very steep on sides, so $\alpha \frac{\partial f}{\partial x}(x^{[k]})$ could be large. This could cause $x^{[1]}$ to overshoot the minimum. This could cause (1) instability, (2) limit cycles, (3) extremely slow and oscillatory convergence

Gradient Descent

Lesson: Don't make α too large.

However if α is too small, then convergence will be very slow.

Most implementations adapt the size of α .

Section 2

Gradient Descent: Multivariable Case

Gradient Descent: Multivariable Case

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a multivariable function.

Example

If $x \in \mathbb{R}^n$ then $f(x) = x_1^2 + x_2^2 + \cdots + x_n^2$ maps $\mathbb{R}^n \rightarrow \mathbb{R}$.

The gradient of a multivariable function is

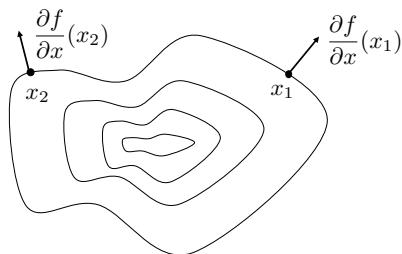
$$\frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

and maps $\mathbb{R}^n \rightarrow \mathbb{R}^n$.

Gradient Descent: Multivariable Case

Example

$$\text{If } f(x) = x_1^2 + \cdots + x_n^2 \text{ then } \frac{\partial f}{\partial x} = \begin{pmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_n \end{pmatrix}$$



The gradient points perpendicular to the level curves of f .

Gradient Descent: Multivariable Case

Theorem (Moon Theorem 14.5)

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a differentiable function in some open set D . The gradient $\frac{\partial f}{\partial \mathbf{x}}(\mathbf{x})$ points in the direction of the maximum increase of f at the point \mathbf{x} .

Gradient Descent: Multivariable Case

Proof.

Expand $f(x + \lambda y)$ in a Taylor series as

$$f(x + \lambda y) = f(x) + \lambda \frac{\partial f^T}{\partial x}(x)y + \text{Higher Order Terms (HOT)}$$

where HOT. are $O(\lambda^2)$, i.e.,

$$\lim_{\lambda \rightarrow 0} \frac{H.O.T.}{\lambda} = 0.$$

We would like to find y that maximizes $f(x + \lambda y)$ as $\lambda \rightarrow 0$.

By Cauchy-Schwartz, $\frac{\partial f^T}{\partial x} y$ is maximized when $y = \frac{\partial f}{\partial x}$. □

Gradient Descent: Multivariable Case

For multivariable functions, the gradient descent formula is

$$x^{[k+1]} = x^{[k]} - \alpha_k \frac{\partial f}{\partial x}(x^{[k]})$$

Again, the selection of the step size is very important. If α_k is too small convergence will be slow.

If α_k is too large, algorithm could be unstable.

How to pick the right α ?

Gradient Descent: Multivariable Case

Locally around a min or max, every smooth function can be approximated by a quadratic (Taylor series).

We can gain insight about the selection of α by studying quadratic functions.

Let $f(x) = x^T R x - 2b^T x$ where $x \in \mathbb{R}^m, b \in \mathbb{R}^m, R = R^T > 0$.

Taking the gradient we get

$$\frac{\partial f}{\partial x} = 2Rx - 2b.$$

Gradient Descent: Multivariable Case

So the gradient descent algorithm is

$$x^{[k+1]} = x^{[k]} - 2\alpha(Rx^{[k]} - b).$$

Let x^* satisfy $Rx^* = b$ then

$$x^{[k+1]} - x^* = x^{[k]} - x^* - 2\alpha(Rx^{[k]} - Rx^*)$$

Define $y^{[k]} = x^{[k]} - x^*$ and $\mu = 2\alpha$, then

$$\begin{aligned}y^{[k+1]} &= y^{[k]} - \mu Ry^{[k]} \\&= (I - \mu R)y^{[k]} \\ \implies y^{[k]} &= (I - \mu R)^k y^{[0]}.\end{aligned}$$

Gradient Descent: Multivariable Case

Since R is symmetric positive definite

$$R = Q\Lambda Q^T$$

where Q -orthogonal. Therefore,

$$\begin{aligned}y^{[k]} &= (QQ^T - \mu Q\Lambda Q^T)^k y^{[0]} \\ &= Q(I - \mu\Lambda)^k Q^T y^{[0]}\end{aligned}$$

Letting $z = Q^T y$,

$$z^{[k]} = (I - \mu\Lambda)^k z^{[0]} \tag{1}$$

$$\implies z_i^{[k]} = (1 - \mu\lambda_i)^k z_i^{[0]} \tag{2}$$

which converges if $|1 - \mu\lambda_i| < 1$, $i = 1, \dots, m$.

Gradient Descent: Multivariable Case

Therefore, convergence happens when

$$\begin{aligned} -1 &< 1 - \mu\lambda_i < 1 \\ \iff -2 &< -\mu\lambda_i < 0 \\ \iff 0 &< \mu\lambda_i < 2 \\ \iff 0 &< \mu < \frac{2}{\lambda_i} \end{aligned}$$

Recall that $\lambda_i > 0$ when R is positive definite, so if

$$0 < \alpha < \frac{1}{\lambda_{\max}(R)}$$

then steepest descent converges for quadratic functions.

Gradient Descent: Multivariable Case

Note that the convergence along each eigenaxis is determined by $\frac{1}{\lambda_i}$.

Therefore if R is ill-conditioned, i.e., $\frac{\lambda_{\max}}{\lambda_{\min}}$ is large, then convergence for gradient descent will be much slower along some axes than others.