



Lecture 9

Math Foundations Team



BITS Pilani

Pilani | Dubai | Goa | Hyderabad



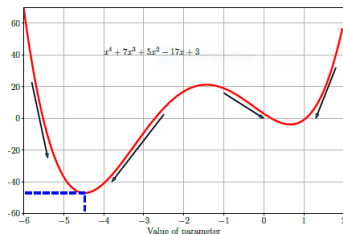
- ▶ We will look at continuous optimization concepts in this lecture.
- ▶ There are two main branches of continuous optimization - constrained and unconstrained optimization.
- ▶ We seek the minimum of an objective function which we assume is differentiable.
- ▶ This is like finding the valleys of the objective function, and since the objective function is differentiable, the gradient tells us the direction to move to get the maximum increase in the objective function

Consider the data in the given table

x	y
1	3.1
2	4.9
3	7.3
4	9.1

- ▶ Easiest model of y one can think of is $y = ax + b$. Let $\hat{y} = ax + b$ be the y predicted.
- ▶ Our aim is to find a and b such that the difference between actual y and the predicted y is minimum. The loss function can be defined as $L(a, b) = \sum_{i=1}^4 (y_i - (ax_i + b))^2$. Then the problem will be to find a and b such that $L(a, b)$ is minimized which is an optimization problem.

- ▶ We move in the direction of the negative gradient to decrease the objective function.
- ▶ We move until we encounter a point at which the gradient is zero.



- ▶ Let $l(x) = x^4 + 7x^3 + 5x^2 - 17x + 3$.
- ▶ The gradient is $\frac{dl(x)}{dx} = 4x^3 + 21x^2 + 10x - 17$.
- ▶ Setting the gradient to zero identifies points corresponding to a local minimum or local maximum - there are three such points since this is a cubic equation.
- ▶ The second derivative is $12x^2 + 42x + 10$



- ▶ For low-order polynomials we can solve the equations analytically and find points at which the gradient is zero.
- ▶ Consider the problem of solving for the minimum of a real-valued function $\min_{\mathbf{x}} f(\mathbf{x})$ where $f : R^d \rightarrow R$ is an objective loss function.
- ▶ We assume our function f is differentiable but that the minimum cannot be found analytically in closed form.
- ▶ The main idea of gradient-descent is to take a step from the current point of magnitude proportional to the negative gradient of the function at the current point.



- ▶ if $\mathbf{x}_1 = \mathbf{x}_0 - \alpha((\nabla f)(\mathbf{x}_0))^T$ for a small step-size $\alpha > 0$ then $f(\mathbf{x}_1) \leq f(\mathbf{x}_0)$.
- ▶ start at some initial point \mathbf{x}_0 and then iterate according to $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha_i((\nabla f)(\mathbf{x}_i))^T$
- ▶ For a suitable step-size α_i , the sequence of points $f(\mathbf{x}_0) \geq f(\mathbf{x}_1) \geq \dots$ converges to some local minimum.
- ▶ α is also called as the learning-rate

Example

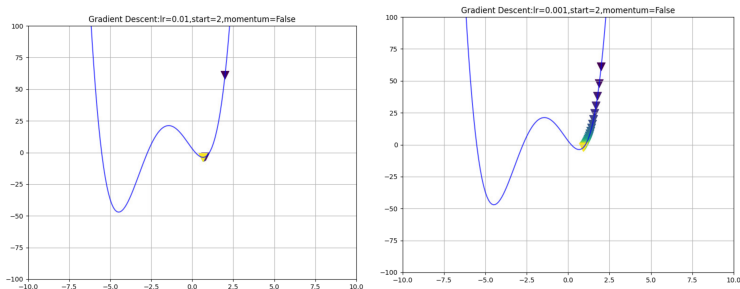


Figure: **left:** with a learning rate of 0.01, local minimum is reached within a couple of steps. **right:** When learning rate is reduced to 0.001, we need relatively more steps to reach the local minimum

Example

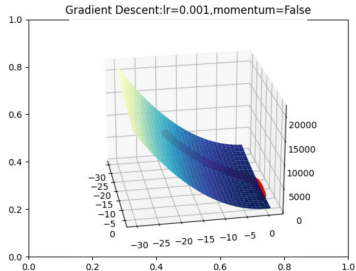
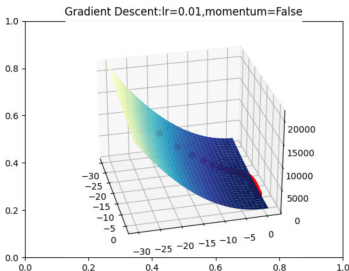


Figure: **left:** with a learning rate of 0.01, minimum is reached. **right:** When learning rate is reduced to 0.001, we need relatively more steps to reach the minimum



- ▶ Consider a machine learning problem consisting of loss functions incurred at N data points.
- ▶ Let the loss function at the i th data point be $L_i(\theta)$
- ▶ Total loss $L(\theta) = \sum_{i=1}^N L_i(\theta)$.
- ▶ Here θ is the parameter vector of interest
- ▶ The standard gradient descent procedure is a batch optimization method $\theta_{i+1} = \theta_i - \alpha_i \sum_{i=1}^N \nabla L_i(\theta_i)^T$.



- ▶ Let S be a subset of the indices $\{1, 2, \dots, N\}$.
- ▶ The set S of data points can be treated as a sample and a sample-centric objective function can be constructed as follows:

$$L_S(\theta) = \sum_{i \in S} L_i(\theta)$$

- ▶ The update equation in case of mini-batch stochastic gradient descent can be written as

$$\theta_{i+1} = \theta_i - \alpha_i \sum_{i \in S} \nabla L_i(\theta_i)^T$$

- ▶ This approach is referred to as mini-batch stochastic gradient.



- ▶ In the extreme case S can contain only one index chosen at random, and the approach is then called as stochastic gradient descent.
- ▶ The key idea in stochastic gradient descent is that the gradient of the sample-specific objective function is an excellent approximation of the true gradient.
- ▶ We can show that when the learning rate decreases at a suitable rate and some mild assumptions can be made, stochastic gradient descent almost surely converges to a local minimum.

Example

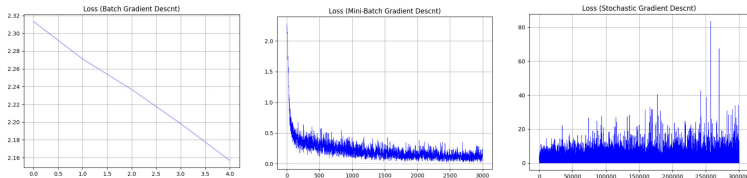


Figure: loss vs num_updates (num epochs:5, dataset:MNIST, layers: lin-relu-lin-relu-lin-relu, loss:crossEntropy, opt:Adam) **left:** Batch Gradient Descent. Entire data is used for every update (thus, 5 epochs results in 5 updates). **right:** Stochastic Gradient Descent. Every update is done based on single sample only. **centre:** Minibatch Gradient Descent. Every update is done using a batch of 100 samples.

Example

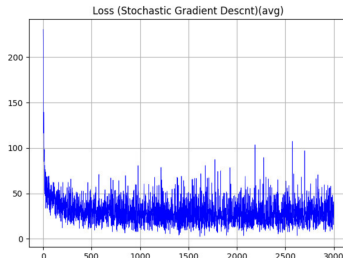


Figure: Though the loss update is done for every sample in SGD, this plot shows the loss averaged over 100 such updates.

Learning rate Algorithm 1 : Decay



- ▶ How are we to decide the value of the learning rate?
- ▶ What happens if we choose a large value for the learning rate and let it be constant? In this case, the algorithm might come close to the optimal answer in the very first iteration but it will then oscillate around the optimal point.
- ▶ What happens if we choose a small value for the learning rate and let it be constant? In this case, it will take a very long time for the algorithm to converge to the optimal point.

Learning rate Algorithm 1 : Decay



- ▶ Choose a variable learning-rate - large initially but decaying with time.
- ▶ This will enable the algorithm to make large strides towards the optimal point and then slowly converge.
- ▶ With a learning-rate dependent on time, the update step becomes $\theta_{t+1} = \theta_t - \alpha_t \nabla L$.

Learning rate Algorithm 1 : Decay



- ▶ The two most common decay functions are exponential decay and inverse decay, expressed mathematically as follows:

exponential decay: $\alpha_t = \alpha_0 e^{-kt}$

inverse decay: $\alpha_t = \frac{\alpha_0}{1 + kt}$

- ▶ In both of the above functions k controls the rate of decay.
- ▶ Another kind of decay function is the step decay where we reduce the learning rate by a constant factor every few steps of gradient descent.

Learning rate Algorithm 2 : Line search



- ▶ Line search uses the optimum step-size directly in order to provide the best improvement.
- ▶ It is rarely used in vanilla gradient descent because of its computational expense, but is helpful in some specialized variations of gradient descent.
- ▶ Let $L(\theta)$ be the function being optimized, and let $\mathbf{d}_t = -\nabla L(\theta_t)$.
- ▶ The update step is $\theta_{t+1} = \theta_t + \alpha_t \mathbf{d}_t$.
- ▶ In line search the learning rate α_t is chosen at the t^{th} step so as to minimize the value of the objective function at θ_{t+1} .
- ▶ Therefore the step-size α_t is computed as $\alpha_t = \min_{\alpha} L(\theta_t + \alpha \mathbf{d}_t)$.

Iteration i of gradient descent $L(x_1, x_2) = x_1^2 + 3x_2^2$

1. At iteration i , $\nabla L(x_1, x_2) = [2x_1, 6x_2]$
2. $H(\alpha) = L(\mathbf{x} - \alpha((\nabla L)(\mathbf{x}))^T) = (1 - 2\alpha)^2 x_1^2 + 3(1 - 6\alpha)^2 x_2^2$.
3. $H'(\alpha) = -4(1 - 2\alpha)x_1^2 - 36(1 - 6\alpha)x_2^2 = 0$.
4. Step Size : $\alpha = \frac{x_1^2 + 9x_2^2}{2x_1^2 + 54x_2^2}$.



- ▶ How do we perform the optimization $\min_{\alpha} L(\theta_t + \alpha \mathbf{d}_t)$?
- ▶ An important property that we exploit of typical line-search settings is that the objective function is a unimodal function of α .
- ▶ This is especially true if we do not use the original objective function but quadratic or convex approximations of it.
- ▶ The first step in optimization is to identify a range $[0, \alpha_{\max}]$ in which to perform the search for the optimum α .
- ▶ We can sweep evaluate the objective function values at geometrically increasing values of α . It is then possible to narrow the search interval by using binary-search or golden-section search method



- ▶ Initialize the search interval $[a, b] = [0, \alpha_{\max}]$.
- ▶ Evaluate the objective function at $\frac{a+b}{2}$ and $\frac{a+b+\epsilon}{2}$
- ▶ Find out whether the function is increasing or decreasing at $\frac{a+b}{2}$. Here ϵ is a small value such as 10^{-6} .
- ▶ If the objective function is found to be increasing at $\frac{a+b}{2}$, we narrow the interval to $[a, \frac{a+b+\epsilon}{2}]$ and continue the search.
- ▶ Otherwise we narrow the interval to $[\frac{a+b}{2}, b]$ and continue the search.



- ▶ Initialize the search interval to $[a, b] = [0, \alpha_{\max}]$.
- ▶ we use the fact that for any mid-samples m_1, m_2 in the region $[a, b]$ where $a < m_1 < m_2 < b$, at least one of the intervals $[a, m_1]$ or $[m_2, b]$ can be dropped. Sometimes we can go so far as to drop $[a, m_2]$ and $[m_1, b]$.
- ▶ When $\alpha = a$ yields the minimum for the objective function, i.e $H(\alpha)$, we can drop the interval $(m_1, b]$.
- ▶ Similarly when $\alpha = b$ yields the minimum for $H(\alpha)$ we can drop the interval $[a, m_2)$. When $\alpha = m_1$ is the value at which the minimum is achieved we can drop $(m_2, b]$.
- ▶ When $\alpha = m_2$ is the value at which the minimum is achieved we can drop $[a, m_1)$.



- ▶ The new bounds on the search interval $[a, b]$ are reset based on the exclusions mentioned in the previous slide. At the end of the process we are left with an interval containing 0 or 1 evaluated point.
- ▶ If we have an interval containing no evaluated point, we select a random point $\alpha = p$ in the reset interval $[a, b]$, and then another point q in the larger of the intervals $[a, p]$ and $[p, b]$.
- ▶ On the other hand if we are left with an interval $[a, b]$ containing a single evaluated point $\alpha = p$, then we select $\alpha = q$ in the larger of the intervals $[a, p]$ and $[p, b]$.
- ▶ This yields another four points on which to continue the golden-section search. We continue until we achieve the desired accuracy.

When do we use line search?



- ▶ The line-search method can be shown to converge to a local optimum, but it is computationally expensive. For this reason, it is rarely used in vanilla gradient descent.
- ▶ Some methods like Newton's method, however, require exact line search.
- ▶ One advantage of using exact line search is that fewer steps are needed to achieve convergence to a local optimum. This might more than compensate for the computational expense of individual steps.