

# Large Scale Optimization for Machine Learning: Lecture 4

Lecturer: Meisam Razaviyayn    Scribe: Sina Baharlouei

Sep 5, 2017

## 1 Recap and Agenda

In the previous lectures, we defined concepts of local and global optimum, convex sets and convex functions. Also, we discussed the necessary conditions for a point to be optimum, the relationship between local and global optimum for a convex function and some important properties of convex sets and convex functions.

In this lecture, we mainly focus on the "Iterative Descent" methods which are a family of Algorithms for finding local optimum of a convex function.

## 2 Overview

As it is mentioned in the previous lecture, if  $\nabla f(x) = 0$ , then  $x$  could be considered as a candidate local minimum point. Now suppose we have a point in which  $\nabla f(x) \neq 0$ . It cannot be a local optimum point, but how can we find another point like  $x'$  so that  $f(x') < f(x)$  using gradient of function on the current point? Intuitively, as we move in the opposite direction of gradient, we get closer to the local minimum point (Gradient vector of a function at a certain point is perpendicular to the tangent line at that point). The following proposition gives the mathematical equivalence of this intuition:

**Proposition 1.** *If  $\nabla f(x)^T d < 0$  then  $\exists \delta$  such that  $f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x}) \forall \alpha \in (0, \delta)$ .*

*Proof.* By the mean value theorem, we can write  $f(\mathbf{x} + \alpha \mathbf{d}) = f(\mathbf{x}) + \alpha \mathbf{d}^T \nabla f(\boldsymbol{\xi})$  where  $\boldsymbol{\xi} = \gamma \mathbf{x} + (1 - \gamma)(\mathbf{x} + \alpha \mathbf{d})$  for some  $\gamma \in (0, 1)$ . We can

choose  $\alpha$  small enough that  $\boldsymbol{\xi}$  is near  $\boldsymbol{x}$  and  $\boldsymbol{d}^T \nabla f(\boldsymbol{\xi}) < 0$ . This gives us  $f(\boldsymbol{x} + \alpha \boldsymbol{d}) < f(\boldsymbol{x})$ .  $\square$

### 3 Iterative Descent Algorithms

Generally speaking, Iterative Descent Algorithms consist of two main components: Direction sequence, Step-size sequence. By choosing these two sequences, an iterative descent algorithm would be obtained. Also these algorithms have an update rule as follow:

$$x^{r+1} \leftarrow x^r + \alpha^r d^r \quad \forall r \in \{0, 1, \dots\} \tag{1}$$

So for each Iterative Algorithm,  $\{\alpha^r\}$  which is the step-size sequence and  $\{d^r\}$  which is the direction sequence, should be proposed.

#### 3.1 Choices of direction

Each direction sequence with elements each of which has an angle more than  $\frac{\pi}{2}$  with the gradient vector on that point, could be considered as a candidate sequence for direction. In this section we introduce some popular directions, widely used in iterative descent Algorithms:

##### 3.1.1 Gradient Descent

A natural way of selecting direction sequence is to choose direction exactly in the opposite direction of gradient vector in each iteration. This choice gives us Gradient/Steepest Descent Algorithm:

$$d^r = -\nabla f(x^r) \tag{2}$$

##### 3.1.2 Newton Direction

Based on the second order Taylor expansion, we can find a more accurate direction for converging, rather than gradient descent's direction. At  $x = x'$ ,  $f(x)$  can be approximated by:

$$f(x) \approx g(x) = f(x') + \nabla f(x')^T (x - x') + \frac{1}{2} (x - x')^T \nabla^2 f(x') (x - x')$$

It is minimized by solving  $\nabla g(x) = 0$ :

$$\nabla g(x) = \nabla f(x') + \nabla^2 f(x')(x - x') = 0$$

which yields to:

$$x = x' - (\nabla^2 f(x'))^{-1} \nabla f(x')$$

So for estimating  $x^{r+1}$  we can use the following formula for direction:

$$d^r = (\nabla^2 f(x^r))^{-1} \nabla f(x^r) \quad (3)$$

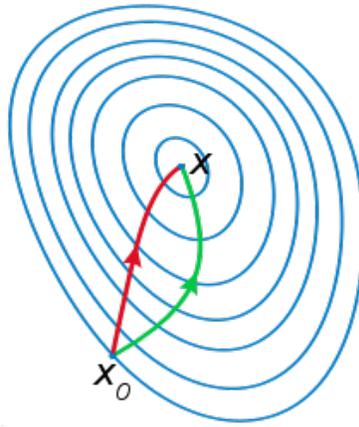


Figure 1: Newton Optimization(Red) versus Gradient Descent(Green)

### 3.1.3 Diagonally scaled gradient descent

Though using Newton method rather than gradient descent, reduces the number of iterations to converge, the cost of computing Newton direction is a drawback; since computation of the inverse of Hessian matrix is a time-consuming procedure. Instead, an approximation of Newton method could be used to reduce the computational complexity. In Diagonally scaled gradient descent, we use the diagonal matrices as the direction sequence.

$$d^r = -D^r \nabla f(x^r) \quad \text{for some } D^r > 0 \quad (4)$$

For instance, we can use the diagonal elements of Newton direction as the direction sequence of Algorithm. As we said, this choice reduces the

computation complexity of Newton method.

$$D^r = \text{diag}(\nabla^2 f(x^r))^{-1} \quad (5)$$

## 3.2 Choices of step-size

In the last section, some well-known directions for iterative descent algorithms has been proposed. In this part, we introduce some popular choices of step-size in iterative descent methods.

### 3.2.1 Constant step-size

A simple way for selecting step-size is to choose a constant value for step-size.

$$\alpha^r = \alpha \quad \forall r \in \{1, 2, \dots\} \quad (6)$$

This method is widely used, because of its simplicity. Picking a large  $\alpha$  might prevent Algorithm to converge; however, choosing a small  $\alpha$  reduces the rate of convergence.

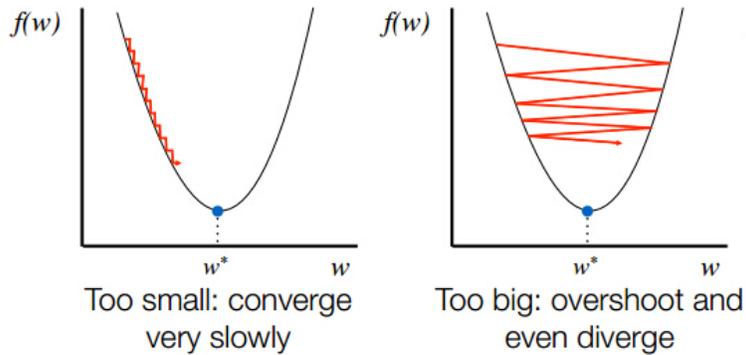


Figure 2: Constant step-size

### 3.2.2 Exact minimization of step-size

One intuitive approach for the choice of step-size is to find the best step-size in each iteration as follow:

$$\alpha^r \in \arg \min_{\alpha \geq 0} f(x^r + \alpha d^r) \quad (7)$$

**Remark.** *This method chooses the best step-size in each iteration locally. So it does not necessarily guarantee to achieve best sequence choice globally.*

### 3.2.3 Limited minimization of step-size

This approach is similar to the exact minimization of step-size. The only difference is that  $\alpha$  can be chosen from a bounded set as follow:

$$\alpha^r \in \arg \min_{\alpha \in (0, \bar{\alpha}]} f(x^r + \alpha d^r) \quad (8)$$

This extra condition prevents over-optimization of the algorithm.

### 3.2.4 Diminishing method

Another way for choosing step-size sequence is to make step-size smaller in each step-size. This method guarantees convergence of the algorithm if it satisfies following conditions:

- $\lim_{r \rightarrow \infty} a^r = 0$
- $\sum_r a^r = \infty$

For example  $\{a_r\} = \frac{1}{r}$  is an acceptable choice. But  $\{a_r\} = \frac{1}{r^2}$  or  $\{a_r\} = \frac{1}{2^r}$  are not acceptable. Since:

$$\sum_r \frac{1}{r^2} = \frac{\pi^2}{6}$$
$$\sum_r \frac{1}{2^r} = 1$$

One important question might be asked is why the second condition is necessary? This is because without having the second condition the set of points that algorithm can be reached would be bounded. And so it might not converge to optimal point.

### 3.2.5 Back tracking (Armijo)

According to the first order Taylor expansion:  $f(\mathbf{x} + \alpha \mathbf{d}) = f(\mathbf{x}) + \alpha \mathbf{d}^T \nabla f(\mathbf{x}) + o(\alpha)$  which implies

$$\underbrace{-\alpha \mathbf{d}^T \nabla f(\mathbf{x})}_{\text{predicted decrease}} \approx \underbrace{f(\mathbf{x}) - f(\mathbf{x} + \alpha \mathbf{d})}_{\text{actual decrease}}.$$

The intuition behind Armijo method is to find maximum step-size among  $\bar{\alpha}$ ,  $\bar{\alpha}\beta$ ,  $\bar{\alpha}\beta^2 \dots$  such that the actual decrease of function will be more than a coefficient of predicted decrease. As we cannot guarantee that predicted decrease is less than actual decrease, the coefficient ( $\sigma$ ) must be less than 1. Mathematically speaking:

**Definition 1.** Let  $\beta, \sigma \in (0, 1)$  be constant numbers, and  $\bar{\alpha}$  be initial step-size, Armijo back-tracking method sequence would be defined as follow:

$$a^r = \max\{\bar{\alpha}\beta^i | f(x^r) - f(x^r + \bar{\alpha}\beta^i d^r) \geq -\sigma \bar{\alpha}\beta^i \nabla f(x^r)^T d^r, i = 0, 1, \dots\}$$

**Proposition 2.** If  $\nabla f(x^r)^T d^r < 0$  then  $a^r$  is well-defined.

## 3.3 Convergence Analysis

In this section we introduce different aspects of analyzing the convergence of an iterative descent Algorithm. One important question about an optimization algorithm is that whether it converges to a stationary point at all? Also it is crucial to ask about the rate of convergence of an Algorithm.

### 3.3.1 Asymptotic rate of convergence

This measurement tends to analyze the behavior of Algorithm asymptotically. In other words, how fast the Algorithm converges to an optimum point in infinity? So first an error function would be defined. For example it could be  $e(x) = \|x - x^*\|$  or  $e(x) = f(x) - f(x^*)$ . The asymptotic behavior of function would be determined as follow:

$$\limsup_{r \rightarrow \infty} \frac{e(\mathbf{x}^{r+1})}{e(\mathbf{x}^r)} = \beta$$

There is three possibility for  $\beta$ :

- Sublinear( $\beta = 1$ ):  $e(x^r) = \frac{1}{r} \Rightarrow \limsup_{r \rightarrow \infty} \frac{e(x^{r+1})}{e(x^r)} = 1.$
- Linear( $\beta \in (0, 1)$ ):  $e(x^r) = \left(\frac{1}{2}\right)^r \Rightarrow \limsup_{r \rightarrow \infty} \frac{e(x^{r+1})}{e(x^r)} = \frac{1}{2}$
- Superlinear( $\beta = 0$ ):  $e(x^r) = \left(\frac{1}{2}\right)^{r^2} \Rightarrow \limsup_{r \rightarrow \infty} \frac{e(x^{r+1})}{e(x^r)} = \left(\frac{1}{2}\right)^{2r+1} = 0$

### 3.3.2 Iteration complexity analysis

If an Algorithm converges faster than another one, It does not mean necessarily the number of iterations for first Algorithm is less than the second one (Because this measure only considers the rate of convergence asymptotically). Instead we can analyze the number of iterations required to achieve  $\epsilon$  – *optimal* solution which means  $e(x^r)$  would be less than  $\epsilon$ . Currently, most of the researchers analyze the worst case number of iterations for an Algorithm.

This measurement tends to analyze the behavior of Algorithm

### 3.3.3 Convergence to a stationary point

$f(x^{r+1}) < f(x^r)$  does not necessarily imply  $\{x^r\}$  converging to a local optimal point even if your function is convex or your step size is small. For example, let  $f(x) = x^2$ . If we let  $x^r = (-1)^r \left(1 + \frac{1}{r}\right)$ , then  $f(x^r) = \left(1 + \frac{1}{r}\right)^2$ . The objective value decreases as you iterate but  $\{x^r\}$  does not converge to a stationary solution.

### 3.3.4 Gradient Related condition

The gradient related condition is satisfied if for any sub-sequence of  $\{x^r\}_{r \in \kappa}$  converging to a non-stationary point, the corresponding sub-sequence of directions  $\{d^r\}_{r \in \kappa}$  is bounded and  $\limsup_{r \rightarrow \infty, r \in \kappa} \nabla f(x^r)^T d^r < 0$  For example:

$$d^r = -D^r \nabla f(x^r) \text{ with } \bar{\gamma}I \geq D^r \geq \gamma I > 0, \forall r$$

The second condition in the example, shows the elements of the diagonal matrix are between  $\gamma$  and  $\bar{\gamma}$ .

- **Descent Lemma.** Suppose we have a smooth function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$ . Then for any  $\mathbf{h}$ , we have  $f(\mathbf{x} + \mathbf{h}) \leq f(\mathbf{x}) + \mathbf{h}^T \nabla f(\mathbf{x}) + \frac{L}{2}\|\mathbf{h}\|^2$ .

Proof. Let  $g(t) = f(\mathbf{x} + t\mathbf{h})$ . Since  $g(1) - g(0) = \int_0^1 g'(t) dt$ , we have

$$\begin{aligned}
f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) &= \int_0^1 \mathbf{h}^T \nabla f(\mathbf{x} + t\mathbf{h}) dt \\
&= \int_0^1 \mathbf{h}^T \nabla f(\mathbf{x}) dt + \int_0^1 \mathbf{h}^T (\nabla f(\mathbf{x} + t\mathbf{h}) - \nabla f(\mathbf{x})) dt \\
&\leq \mathbf{h}^T \nabla f(\mathbf{x}) + \int_0^1 \|\mathbf{h}\| \cdot \|\nabla f(\mathbf{x} + t\mathbf{h}) - \nabla f(\mathbf{x})\| dt \\
&\leq \mathbf{h}^T \nabla f(\mathbf{x}) + \int_0^1 tL\|\mathbf{h}\|^2 dt \\
&= \mathbf{h}^T \nabla f(\mathbf{x}) + \frac{L}{2}\|\mathbf{h}\|^2.
\end{aligned}$$

- **Theorem**

Assume  $\mathbf{x}^{r+1} \leftarrow \mathbf{x}^r + \alpha^r \mathbf{d}^r$ ,  $\mathbf{d}^r$  is gradient related, function  $f$  has Lipschitz gradient and one of the following step size rules hold:

(a) Diminishing:  $\alpha^r \rightarrow 0$  and  $\sum_r \alpha^r = \infty$

(b) Armijo

(c) Small enough:  $0 < \epsilon \leq \alpha^r \leq \frac{(2-\epsilon)\|\nabla f(\mathbf{x}^r)^T \mathbf{d}^r\|}{L\|\mathbf{d}^r\|^2}$

Then every limit point of the iterates is a stationary point, i.e.

$$\text{if } \{\mathbf{x}^r\}_{r \in \kappa}, \text{ then } \nabla f(\bar{\mathbf{x}}) = 0.$$

- Proof of part (a)

Proof is by contradiction. Assume the contrary that  $\{\mathbf{x}^r\}_{r \in \kappa} \rightarrow \bar{\mathbf{x}}$  but  $\bar{\mathbf{x}}$  is not stationary. Therefore, by the gradient related condition, we should have  $\nabla f(\mathbf{x}^r)^T \mathbf{d}^r \leq C < 0$  for large enough  $r$ . Furthermore, after

some point, the step-size  $\alpha^r$  is small enough and by the descent lemma, we have  $f(\mathbf{x}^{r+1}) - f(\mathbf{x}^r) \leq \frac{\epsilon}{2} \alpha^r \nabla f(\mathbf{x}^r)^T \mathbf{d}^r \leq \frac{\epsilon}{2} \alpha^r C \leq 0$ . Summing up all inequalities, we have that

$$\begin{aligned} \sum_{r=0}^t f(\mathbf{x}^{r+1}) - f(\mathbf{x}^r) &\leq \sum_{r=0}^t \frac{\epsilon}{2} \alpha^r \nabla f(\mathbf{x}^r)^T \mathbf{d}^r \leq \sum_{r=0}^t \frac{\epsilon}{2} \alpha^r C \leq 0. \\ \Rightarrow f(\mathbf{x}^{t+1}) - f(\mathbf{x}^0) &\leq \frac{\epsilon}{2} C \sum_{r=0}^t \alpha^r \leq 0. \end{aligned}$$

By letting  $t \rightarrow \infty$  we have,

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^0) \leq \frac{\epsilon}{2} C \sum_{r=0}^{\infty} \alpha^r \leq 0$$

from which we conclude  $\sum_{r=0}^{\infty} \alpha^r$  is upper bounded by a constant which contradicts the assumption that  $\sum_{r=0}^{\infty} \alpha^r = \infty$ .

- Proof of part (b)

Proof by contradiction: assume the contrary that  $\{x^r\} \rightarrow \bar{x}$  and  $\bar{x}$  is not a stationary point. Let us assume that  $\alpha^r = \beta^{i_r}$  with  $\beta$  is as defined in the definition of Armijo rule. If  $i_r = 0, \forall r \geq r_0$  then the proof is the same as the proof in part (c). Therefore, let us assume that  $i_r \geq 1$  for a subsequence of iterates. By restricting to a subsequence if necessary, we can write

$$f(x^r) - f(x^r + \alpha^r \mathbf{d}^r) \geq -\sigma \alpha^r \nabla f(x^r)^T \mathbf{d}^r \geq 0, \quad (9)$$

and

$$f(x^r) - f(x^r + \frac{\alpha^r}{\beta} \mathbf{d}^r) < -\sigma \frac{\alpha^r}{\beta} \nabla f(x^r)^T \mathbf{d}^r. \quad (10)$$

Defining  $p^r = \frac{\mathbf{d}^r}{\|\mathbf{d}^r\|}$  and  $\bar{\alpha}^r = \alpha^r \frac{\|\mathbf{d}^r\|}{\beta}$ , (10) can be rewritten as

$$\frac{f(x^r) - f(x^r + \bar{\alpha}^r p^r)}{\bar{\alpha}^r} < -\sigma \nabla f(x^r)^T p^r. \quad (11)$$

Since  $\{p^r\}$  belongs to a compact set (the surface of unit ball), it has a limit point  $\bar{p}$ . On the other hand, letting  $r \rightarrow \infty$  in (9) implies that  $\alpha^r \mathbf{d}^r \rightarrow 0$ . Taking the limit  $r \rightarrow \infty$  in (11) leads to

$$-\nabla f(\bar{x})^T \bar{p} \leq -\sigma \nabla f(\bar{x})^T \bar{p}.$$

Since  $\sigma < 1$ , we conclude that  $\nabla f(\bar{x})^T \bar{p} \geq 0$  which contradicts the gradient related condition.

- Proof of Part (c)

Proof by contradiction: We have  $f(\mathbf{x}^r + \alpha^r \mathbf{d}^r) - f(\mathbf{x}^r) \leq \alpha^r \nabla f(\mathbf{x}^r)^T \mathbf{d}^r + \frac{L}{2} (\alpha^r)^2 \|\mathbf{d}^r\|^2$  from Descent Lemma. Suppose  $\{\mathbf{x}^r\}_{r \in \kappa} \rightarrow \bar{\mathbf{x}}$  but  $\bar{\mathbf{x}}$  is not stationary.

$$\begin{aligned} f(\mathbf{x}^r + \alpha^r \mathbf{d}^r) - f(\mathbf{x}^r) &\leq \alpha^r \nabla f(\mathbf{x}^r)^T \mathbf{d}^r + \frac{L}{2} \alpha^r \|\mathbf{d}^r\|^2 \frac{(2 - \epsilon) |\nabla f(\mathbf{x}^r)^T \mathbf{d}^r|}{L \|\mathbf{d}^r\|^2} \\ &= \alpha^r \nabla f(\mathbf{x}^r)^T \mathbf{d}^r \left(1 - \frac{2 - \epsilon}{2}\right) \\ &= \frac{\epsilon}{2} \alpha^r \nabla f(\mathbf{x}^r)^T \mathbf{d}^r \leq 0 \end{aligned} \tag{12}$$

which implies that  $f(\mathbf{x}^{r+1}) \leq f(\mathbf{x}^r)$ . If  $\mathbf{x}^r \rightarrow \bar{\mathbf{x}}$  then  $f(\mathbf{x}^r) \rightarrow f(\bar{\mathbf{x}})$ . Therefore by taking limit as  $r \rightarrow \infty$  in (12), we obtain  $\lim_{r \rightarrow \infty} \nabla f(\mathbf{x}^r)^T \mathbf{d}^r = 0$  which contradicts gradient related condition in slide 7.