

Mariia Aleksandrovyh, Yaffa Atkins, Atoosa Lotfi

1 Exercise 1.3

Use a Taylor series expansion to show that for any descent direction $d(\theta)$ at a point θ of a twice continuously differentiable function, there exists $\epsilon_0 > 0$ such that

$$J(\theta + \epsilon d(\theta)) \leq J(\theta) \text{ for all } 0 \leq \epsilon \leq \epsilon_0.$$

(Answer)

Use Taylor Series expansion to show that for any descent direction $d(\theta)$, there exists some positive value where $J(\theta + \epsilon d(\theta))$ shrinks $J(\theta)$ for at least the first two derivatives.

1st derivative: descent direction $d(\theta) = -\nabla J(\theta)$, since the gradient is negative, the first derivative is negative.
Use hessian when multiple dimensions

2nd: Even if positive, its minimized because there's some tiny epsilon between 0 and 1 where $\epsilon \cdot \text{hessian} > \text{gradient}$

Taylor series: $J(\theta + \epsilon \nabla J(\theta)) = J(\theta) + \epsilon \cdot \nabla^2 J(\theta) + \dots$

1.3

$$J(\theta + \varepsilon d) = J(\theta) + \varepsilon d^T \nabla J(\theta) + \frac{\varepsilon^2}{2} d^T \nabla^2 J(\theta) d$$

$\nabla^2 J(\theta)$ Hessian of $J(\theta)$, symmetric matrix $d \times d$ whose (i,j) entry is $\frac{\partial^2 f}{\partial x_i \partial x_j}$

if $J(\theta): \mathbb{R}^d \rightarrow \mathbb{R}$ such that J is

twice-differentiable and has continuous derivatives in an open ball B around point $\theta \in \mathbb{R}^d$

Then for small enough $\varepsilon \in \mathbb{R}^d$ such that

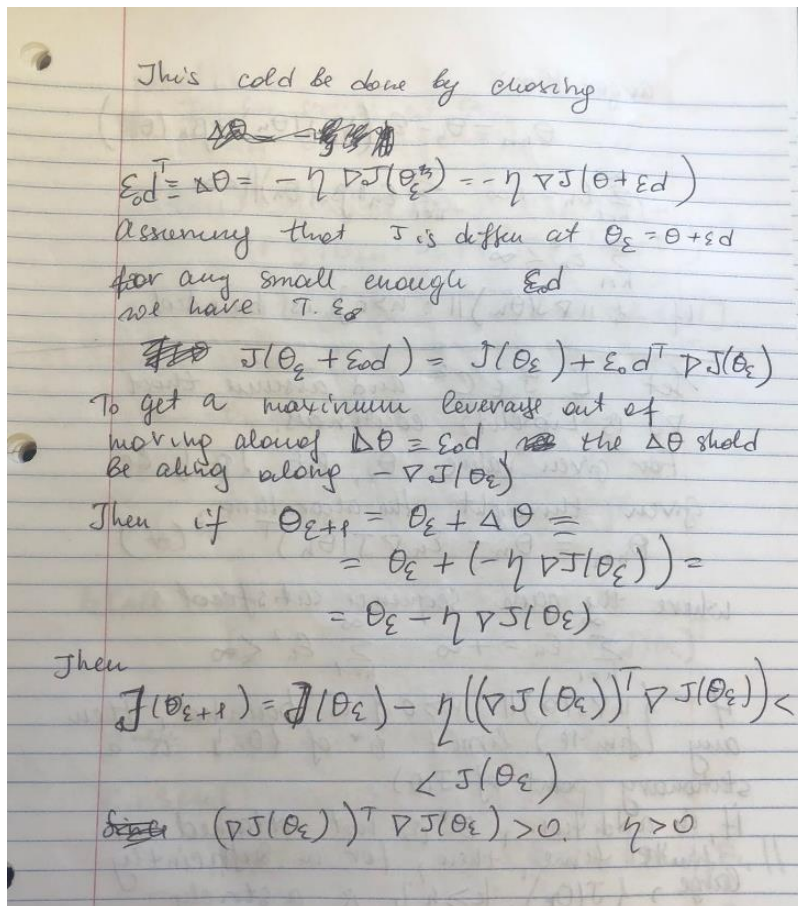
$\theta + \varepsilon d$ is also in \mathbb{R}^d (ball B)

we have following:

$$J(\theta + \varepsilon d) = J(\theta) + \varepsilon d^T \nabla J(\theta) + \frac{\varepsilon^2}{2} d^T \nabla^2 J(\theta) d$$

Then for any step $\varepsilon \geq 0$ $0 \leq \varepsilon \leq \varepsilon_0$ when we at point $\theta' = \theta + \varepsilon d \in \mathbb{R}^d$, gradient decent move has to move in direction $\varepsilon d \in \mathbb{R}^d$ such that

$$J(\theta' + \varepsilon d) < J(\theta')$$



2 Exercise 1.9

Let $J \in C^2$ be such that the gradient is a bounded and Lipschitz continuous function and consider the biased algorithm:

$$\theta_{n+1} = \theta_n - \epsilon_n (\nabla_\theta J(\theta_n) + \beta_n(\theta_n)),$$

where

$$\sum_{n=1}^{\infty} \epsilon_n = +\infty, \quad \sum_{n=1}^{\infty} \epsilon_n \|\beta_n(\theta_n)\| < \infty, \quad \sum_{n=1}^{\infty} \epsilon_n^2 < \infty.$$

If $\{\|\nabla J(\theta_n)\| : n \geq 0\}$ is bounded, then any limit point not achieved in finite time is a stationary point of J .

(Answer)

We use Taylor series to write cost function:

$$J(\theta_{n+1}) = J(\theta_n) + \nabla J(\theta_n)^T (\theta_{n+1} - \theta_n) + \frac{1}{2} (\theta_{n+1} - \theta_n)^T \nabla^2 J(\theta_n) (\theta_{n+1} - \theta_n) \quad (1)$$

Then substitute $(\theta_{n+1} - \theta_n)$

$$\theta_{n+1} = \theta_n - \epsilon_n (\nabla J(\theta_n) + \beta_n(\theta_n)) \Rightarrow (\theta_{n+1} - \theta_n) = -\epsilon_n (\nabla J(\theta_n) + \beta_n(\theta_n)) \quad (2)$$

Then the equation (1) becomes:

$$J(\theta_{n+1}) = J(\theta_n) - \varepsilon_n \|\nabla J(\theta_n)\|^2 - \varepsilon_n \nabla J(\theta_n)^T \beta_n(\theta_n) + \frac{1}{2} \varepsilon_n^2 (\nabla J(\theta_n) + \beta_n(\theta_n))^T \nabla^2 J(\theta_n) (\nabla J(\theta_n) + \beta_n(\theta_n))$$

Now we recall

$$\frac{1}{2} \varepsilon_n^2 (\nabla J(\theta_n) + \beta_n(\theta_n))^T \nabla^2 J(\theta_n) (\nabla J(\theta_n) + \beta_n(\theta_n)) = h_n$$

Then

$$\frac{1}{2} \varepsilon_n^2 (\nabla J(\theta_n) + \beta_n(\theta_n))^T \nabla^2 J(\theta_n) (\nabla J(\theta_n) + \beta_n(\theta_n)) \leq \frac{1}{2} \varepsilon_n^2 L \|\nabla J(\theta_n) + \beta_n(\theta_n)\|^2$$

$$h_n \leq \frac{1}{2} \varepsilon_n^2 L \|\nabla J(\theta_n) + \beta_n(\theta_n)\|^2$$

Let's take a limit

$$\sum_n h_n \leq \sum_n \frac{1}{2} \varepsilon_n^2 L \|\nabla J(\theta_n) + \beta_n(\theta_n)\|^2 \leq \frac{1}{2} L \left(\sum_n \varepsilon_n^2 \|\nabla J(\theta_n)\|^2 + \sum_n \varepsilon_n^2 \|\beta_n\|^2 \right)$$

Where $\sum_n \varepsilon_n^2 \|\nabla J(\theta_n)\|^2 < \infty$ because gradient is bounded, and $\sum_n \varepsilon_n^2 < \infty$.

Also, $\sum_n \varepsilon_n^2 \|\beta_n\| < \infty$

This means that h_n is summable, so we apply Lemma 1.2.

Which give us the next conditions:

$$\text{i.} \quad J(\theta_n) \rightarrow -\infty$$

Or

$$\text{ii.} \quad J(\theta_n) \text{ converges to a finite value and } \sum_n \varepsilon_n (\|\nabla J(\theta_n)\|^2 + \nabla J(\theta_n)^T \beta_n(\theta_n)) < \infty \text{ converges}$$

If $\bar{\theta}$ the limit of the θ_n is not achieved in finite time, then we have the

$$\sum_n \varepsilon_n (\|\nabla J(\theta_n)\|^2 + \nabla J(\theta_n)^T \beta_n(\theta_n)) < \infty \text{ and continuity of } \nabla J(\theta_n) \text{ so that}$$

$$\|\nabla J(\bar{\theta})\| = \lim_{i \rightarrow \infty} \|\nabla J(\theta_{m_i})\| = 0$$

This shows that $\bar{\theta}$ is stationary point.

Let's prove it!

$\bar{\theta}$ is the limit of θ_n . Then using

$$\frac{1}{2} \varepsilon_n^2 (\nabla J(\theta_n) + \beta_n(\theta_n))^T \nabla^2 J(\theta_n) (\nabla J(\theta_n) + \beta_n(\theta_n)) \leq \frac{1}{2} \varepsilon_n^2 L \left\| (\nabla J(\theta_n) + \beta_n(\theta_n)) \right\|^2$$

The next equation

$$\begin{aligned} J(\theta_{i+1}) &= J(\theta_i) - \varepsilon_i \left\| \nabla J(\theta_i) \right\|^2 - \varepsilon_i \nabla J(\theta_i)^T \beta_i(\theta_i) + \\ &+ \frac{1}{2} \varepsilon_i^2 (\nabla J(\theta_i) + \beta_i(\theta_i))^T \nabla^2 J(\theta_i) (\nabla J(\theta_i) + \beta_i(\theta_i)) \end{aligned}$$

Becomes

$$J(\theta_{i+1}) \leq J(\theta_i) - \varepsilon_i \left(\left\| \nabla J(\theta_i) \right\|^2 + \nabla J(\theta_i)^T \beta_i(\theta_i) \right) + \frac{1}{2} \varepsilon_i L \left\| (\nabla J(\theta_i) + \beta_i(\theta_i)) \right\|^2$$

As $i \rightarrow \infty$ and $\varepsilon_i < 1$

$\left(\left\| \nabla J(\theta_i) \right\|^2 + \nabla J(\theta_i)^T \beta_i(\theta_i) \right) + \frac{1}{2} \varepsilon_i L \left\| (\nabla J(\theta_i) + \beta_i(\theta_i)) \right\|^2$ is positive so, i_0 exist such that $\{J(\theta_i) : i \geq i_0\}$ is decreasing towards $J(\bar{\theta})$, where $\bar{\theta}$ is a stationary point.

3 Exercise 1.11

Consider a single machine that can operate one piece at a time, and let the service times of the machine constitute a sequence of iid exponentially distributed random variables. Parts arrive to the machine according to a Poisson process with unit rate. In other words, the time between arrivals of parts is exponentially distributed with mean 1 and that interarrival times are mutually independent. In order for the system to be stable, assume that the expected service time is strictly less than one. Let $C(\theta) = 1/\theta^2$ be the cost of operating the system at service mean θ . Let $P(\theta)$ denote the stationary probability that the queue length is larger than or equal to a threshold b .

(a) Find the solution to the constrained problem:

$$\min C(\theta), \text{ s.t. } P(\theta) \leq \alpha.$$

Interpret the constraint qualifications, the second order condition, and the Lagrange multiplier. Hint: Use the fact that the probability that the stationary queue length equals n is given by $(1 - \theta)\theta^n$, for $n \in \mathbb{N}$.

(Answer)

1.11

Say there's an assembly line. parts come at an average rate of 1 and pass through at a rate faster than 1.

but since both rates are random, there could be pile ups.

$C(\theta)$ = cost of operating the system at service mean θ

Find $\min C(\theta)$ such that $P(\theta) \leq \alpha$

$$P(\theta) = P(\text{queue} = \text{any } n \geq b)$$

$$P(\theta) = \sum_{n=b}^{\infty} p(\theta) = \sum_{n=b}^{\infty} (1-\theta)(\theta^n)$$

sum this formula as n goes from b to infinity

$$P(\theta) = \theta^b - \theta^{\infty} = \theta^b$$

$\theta^{\infty} = 0$ because θ = service time and there's

a constraint that service time < 1 so $\theta < 1$

exponents on θ for $0 < \theta < 1$ get exponentially smaller

so θ^{∞} is basically 0

$$\begin{aligned} \text{so } P(\theta) = \theta^b &\leq \alpha \iff \theta = \alpha^{1/b} \\ &= \theta \leq \alpha^{1/b} \\ &= \frac{\theta}{\theta} \leq \frac{\alpha^{1/b}}{\theta} \\ &= \alpha^{-1/b} \leq \frac{1}{\theta} \\ &= \alpha^{-2/b} \leq \frac{1}{\theta^2} \end{aligned}$$

$$\text{we know } C(\theta) = \frac{1}{\theta^2}$$

$$\text{so } C(\theta) \geq \alpha^{-2/b}$$

$$\min C(\theta) = \alpha^{-2/b}$$

Second order condition:

Cost function:

$$C(\theta) = 1/\theta^2$$

and two constraints:

1. $(\theta^b - \alpha) \leq 0$
2. $\theta - 1 < 0$

Then unconstrained problem:

$$f(\theta) = 1/\theta^2 + \alpha_1 (\theta^b - \alpha)^2 + \alpha_2 (\theta - 1)^2$$

Subject to minimization using second conditions:

$F' = 0 \rightarrow$ stationary point,

$F'' > 0$ or $F'' < 0 \rightarrow$ local minimum or maximum

Lagrange multipliers:

We have two constraints:

$$3. \quad (\theta^b - \alpha) \leq 0$$

$$4. \quad \theta - 1 < 0$$

Then the Lagrangian function:

$$L(\theta; \lambda_1, \lambda_2) = \frac{1}{\theta^2} + \lambda_1(\theta^b - \alpha) + \lambda_2(\theta - 1)$$

From definition 1.11 we have

$$\nabla_{\theta} \mathcal{L}(\theta^*, \lambda^*, \eta^*) = 0$$

$$\nabla_{\lambda} \mathcal{L}(\theta^*, \lambda^*, \eta^*) = g(\theta^*)^{\top} \leq 0, \lambda^* \geq 0, \text{ and } \forall i : \lambda_i^* g_i(\theta^*) = 0$$

$$\nabla_{\eta} \mathcal{L}(\theta^*, \lambda^*, \eta^*) = h(\theta^*)^{\top} = 0;$$

Then we take partial derivative in relatively to each variable and multiplier:

So, we have:

$$\frac{\partial L}{\partial \theta} = -\frac{2}{\theta^3} + \lambda_1(b\theta^{b-1}) + \lambda_2 = 0 \Rightarrow \lambda_2 = \frac{2}{\theta^3} - \lambda_1(b\theta^{b-1}) \quad (1)$$

$$\frac{\partial L}{\partial \lambda_1} = (\theta^b - \alpha) = g_1(\theta) \Rightarrow \lambda_1 g_1 = \lambda_1(\theta^b - \alpha) = 0 \Rightarrow (\theta^b - \alpha) = 0 \text{ or } \lambda_1 = 0 \Rightarrow \left[\theta = \alpha^{1/b} \right]$$

$$\frac{\partial L}{\partial \lambda_2} = (\theta - 1) = g_2 \Rightarrow \lambda_2 g_2 = \lambda_2(\theta - 1) = 0 \Rightarrow [\lambda_2 = 0]$$

$$\text{And from one we have: } \lambda_1 b \theta^{b-1} = \frac{2}{\theta^3} \quad \lambda_1 = \frac{2}{b \theta^{b-2}}$$

Then with the $\alpha = 0.01$ and $b = 10$ $\theta^* = 0.6310$

- (b) Program two different numerical methods to solve this problem for $\alpha = .01$ and $b = 10$. Plot the consecutive values of θ_n and discuss your results, comparing with the theoretical answer in part (a).

(Answer)

[GoogleColab HW1](#)

https://colab.research.google.com/drive/1u-SxxZ8fDp_VrruJm58eVr-tUfcB0usx?usp=sharing

Discuss results

Compare with part a

Gradient Descent Method:

Using theta of .25, our method converges at $y = 0.9426$ $x = -1.0300$ with $\theta_N = -1.03$. Theta starts positive but over time it is minimized by $\alpha \cdot \text{derivative}$, so it slowly gets smaller and smaller.

In part a we found a positive θ less than 1, that the second constrain is in active and just the first constrain over beta is active.

Gradient Descent

```
[8] import numpy as np
import matplotlib.pyplot as plt
from matplotlib.ticker import MaxNLocator
from itertools import product

def C(theta):    #cost function
    return 1/(theta**2)

def cprime(theta):    #gradient cost function
    return -2/(theta**3)

def plotFunc(theta0):
    theta = np.linspace(-5, 7, 100)
    plt.plot(theta, C(theta))
    plt.plot(theta0, C(theta0), 'ro')
    plt.xlabel('Theta')
    plt.ylabel('f(theta)')
    plt.title('Gradient Descent')

def plotPath(xs, ys, theta0):
    plotFunc(theta0)
    plt.plot(xs, ys, linestyle='--', marker='o', color='blue')
    plt.plot(xs[-1], ys[-1], 'ro')

theta0 = 0.25
plotFunc(theta0)
```



```
def GradientDescent(function, cprime, theta0, alpha, beta):
    #setting the variables
    theta_K = theta0
    c_theta_k = C(theta_K)
    primeC_theta_k = -cpime(theta_K)
    x_axis = [theta_K] #store theta
    y_axis = [c_theta_k] #store c(theta)

    #P(theta) is supposed to be less than or equal to alpha, P(theta) = theta^b
    while theta_K**beta <= alpha:
        theta_K = theta_K + alpha * primeC_theta_k
        c_theta_k = C(theta_K)
        primeC_theta_k = -cpime(theta_K)
        x_axis.append(theta_K) #plot theta
        y_axis.append(c_theta_k) #plot c(theta)

    # print results
    if theta_K**beta == alpha:
        print('Gradient descent does not converge.')
    else:
        print('Solution:\n y = {:.4f}\n x = {:.4f}'.format(c_theta_k, theta_K))
        print(theta_K)
    return x_axis, y_axis

xs, ys = GradientDescent(function, cprime, theta0, a, b)
plotPath(xs, ys, theta0)
```

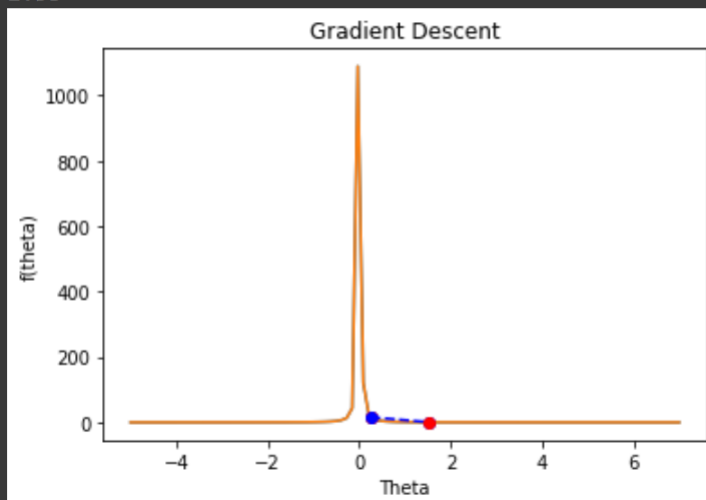


Solution:

y = 0.4272

x = 1.5300

1.53



Also we tried the penalty method but it