

# COM SCI 260B HW 1 Solution

Ashish Kumar Singh (UID:105479019)

April 12, 2022

**Problem 1.** Bound on gradient

**Solution 1a.** One example function is  $f(x, y) = x^2 - y^2$ ,  $\nabla f(x, y) = \begin{bmatrix} 2x \\ -2y \end{bmatrix}$

At  $(x, y) = (0, 0)$  we see that the gradient is zero, but it is not a local maximum or minimum as  $f(0, 0) = 0$  but  $f(\delta, 0) = \delta^2$  and  $f(0, \delta) = -\delta^2$

In the plot below we can see,  $(x, y) = (0, 0)$  is not a maxima or a minima

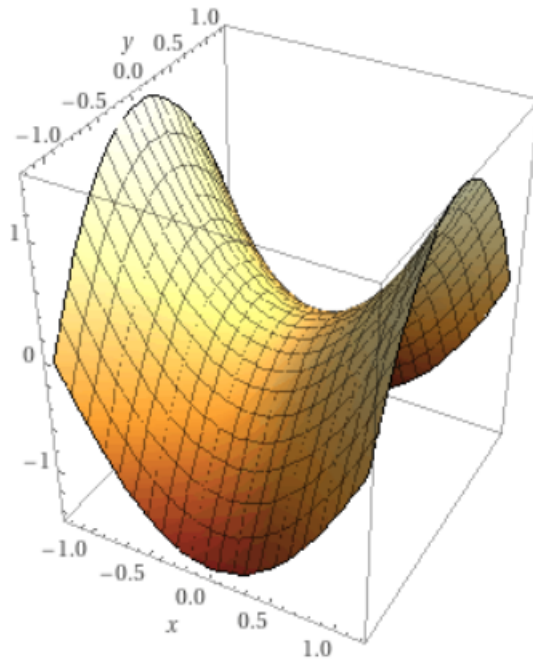


Figure 1: Plot of  $f(x, y) = x^2 - y^2$

**Solution 1b.** From the monotonicity of GD, we have

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2 \quad (1)$$

Lets say GD converges to  $x_*$ , then

$$f(x_*) \leq f(x_k) - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2 \quad (2)$$

Now from the definition of  $\beta$  - *smoothness*, we have

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|_2^2$$

Using  $y = x_k$ ,  $x = x_*$  and  $\nabla f(x_*) = 0$ , we get

$$f(x_k) \leq f(x_*) + \frac{\beta}{2} \|x_k - x_*\|_2^2 \quad (3)$$

Substituting eqn2 in eqn3 we get,

$$f(x_k) \leq f(x_k) - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2 + \frac{\beta}{2} \|x_k - x_*\|_2^2$$

$$\|\nabla f(x_k)\|_2^2 \leq \beta^2 \|x_k - x_*\|_2^2$$

The right handside goes to 0 as we converge, which shows that with GD we can find  $w = x_k$  which has arbitrarily small gradient norm.

Now , lets assume kth iteration of GD is the first iteration when the gradient norm is less then  $\epsilon$ , i.e.  $\|\nabla f(x_k)\|_2 \leq \epsilon$  and  $\|\nabla f(x_i)\|_2 > \|\nabla f(x_k)\|_2$  for every  $i < k$ ,

Writing eqn1 for different values of k we have,

$$\begin{aligned} f(x_1) &\leq f(x_0) - \frac{1}{2\beta} \|\nabla f(x_0)\|_2^2 \\ f(x_2) &\leq f(x_1) - \frac{1}{2\beta} \|\nabla f(x_1)\|_2^2 \\ &\vdots \\ f(x_k) &\leq f(x_{k-1}) - \frac{1}{2\beta} \|\nabla f(x_{k-1})\|_2^2 \end{aligned}$$

Adding the above inequalities, and using  $\|\nabla f(x_i)\|_2 > \|\nabla f(x_k)\|_2$  for every  $i < k$  we get

$$f(x_k) \leq f(x_0) - \frac{k}{2\beta} \|\nabla f(x_k)\|_2^2 \quad (4)$$

Substituting eqn4 in eqn2 we get,

$$f(x_*) \leq f(x_0) - \frac{k}{2\beta} \|\nabla f(x_k)\|_2^2 - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2$$

$$\|\nabla f(x_k)\|_2^2 \leq \frac{2\beta}{k+1}(f(x_0) - f(x_*))$$

For  $\|\nabla f(x_k)\|_2 \leq \epsilon$ , we will have

$$\frac{2\beta}{k+1}(f(x_0) - f(x_*)) \leq \epsilon^2$$

$$k \geq \frac{2\beta}{\epsilon^2}(f(x_0) - f(x_*)) - 1 \tag{5}$$

Thus, with GD we can find points with arbitrarily small gradient norm. The number of steps required to achieve it is bounded by eqn5

**Problem 2.** Convergence rate of GD

**Solution 2.** Given that for  $\alpha$ -convex function  $f$ ,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|_2^2 \quad (6)$$

Also from monotonicity of GD for  $\beta$ -smooth function  $f$ ,

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2 \quad (7)$$

Using  $y = x_*$  and  $x = x_k$  in eqn6 and putting in eqn7,

$$f(x_{k+1}) \leq f(x_*) - \langle \nabla f(x_k), x_* - x_k \rangle - \frac{\alpha}{2} \|x_* - x_k\|_2^2 - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2$$

Using  $\nabla f(x_k) = \beta(x_k - x_{k+1})$

$$\begin{aligned} f(x_{k+1}) - f(x_*) &\leq \frac{1}{2\beta} \left[ 2\beta \langle \beta(x_k - x_{k+1}), x_k - x_* \rangle - \alpha\beta \|x_* - x_k\|_2^2 - \beta^2 \|x_k - x_{k+1}\|_2^2 \right] \\ f(x_{k+1}) - f(x_*) &\leq \frac{\beta}{2} \left[ \|x_k - x_{k+1}\|_2^2 + \|x_* - x_k\|_2^2 - \|x_* - x_{k+1}\|_2^2 - \frac{\alpha}{\beta} \|x_* - x_k\|_2^2 - \|x_k - x_{k+1}\|_2^2 \right] \\ f(x_{k+1}) - f(x_*) &\leq \frac{\beta}{2} \left[ \left(1 - \frac{\alpha}{\beta}\right) \|x_k - x_*\|_2^2 - \|x_{k+1} - x_*\|_2^2 \right] \end{aligned} \quad (8)$$

Note that the left hand side will always be positive because of the monotonicity of GD, so we can write,

$$\begin{aligned} 0 &\leq \frac{\beta}{2} \left[ \left(1 - \frac{\alpha}{\beta}\right) \|x_k - x_*\|_2^2 - \|x_{k+1} - x_*\|_2^2 \right] \\ \|x_{k+1} - x_*\|_2^2 &\leq \left(1 - \frac{\alpha}{\beta}\right) \|x_k - x_*\|_2^2 \end{aligned}$$

For simplicity, let's assume  $\lambda = (1 - \frac{\alpha}{\beta})$ , note that by definition of convexity and smoothness, it is obvious that  $\beta > \alpha$ , hence  $\lambda > 0$ . Writing above eqn for different values of  $k$ , we get

$$\begin{aligned} \|x_1 - x_*\|_2^2 &\leq \lambda \|x_0 - x_*\|_2^2 \\ \|x_2 - x_*\|_2^2 &\leq \lambda \|x_1 - x_*\|_2^2 \\ &\vdots \\ \|x_k - x_*\|_2^2 &\leq \lambda \|x_{k-1} - x_*\|_2^2 \end{aligned}$$

On multiplying the above inequalities, we get

$$\begin{aligned} \|x_k - x_*\|_2^2 &\leq \lambda^k \|x_0 - x_*\|_2^2 \\ \|x_k - x_*\|_2^2 &\leq \left(1 - \frac{\alpha}{\beta}\right)^k \|x_0 - x_*\|_2^2 \end{aligned}$$

**Problem 3.** Differentiable convex function

**Solution 3.** Suppose for function  $f : R^d \rightarrow R$  we have a local minimum at  $x_l$ , hence  $\nabla f(x_l) = 0$

Now, lets assume that it is possible to find a global minimum at  $x_g$ , such that  $f(x_g) < f(x_l)$ .

Using the definition of convexity, for every  $u, v$  we have

$$f(v) \geq f(u) + \langle \nabla f(u), v - u \rangle$$

Using  $v = x_g$  and  $u = x_l$  and  $\nabla f(x_l) = 0$ ,

$$f(x_g) \geq f(x_l) + \langle \nabla f(x_l), x_g - x_l \rangle$$

$$f(x_g) \geq f(x_l)$$

which contradicts our assumption that  $f(x_g) < f(x_l)$ , hence it is not possible to find another point whose function value is lower than  $f(x_l)$ , thus for convex differentiable function every local minimum is the global minimum.

**Problem 4.** GD vs NAGD

**Solution 4a.** Given,

$$L(w) = \frac{1}{n} \sum_{i=1}^n l(y_i, \sigma(\langle w, x_i \rangle))$$

$$L(w) = \frac{-1}{n} \sum_{i=1}^n (y_i \log(\sigma(\langle w, x_i \rangle)) + (1 - y_i) \log(1 - \sigma(\langle w, x_i \rangle)))$$

Differentiating wrt  $w_j$  and using  $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

$$\frac{dL(w)}{dw_j} = \frac{-1}{n} \sum_{i=1}^n x_{ij} (y_i (1 - \sigma(\langle w, x_i \rangle)) - (1 - y_i) \sigma(\langle w, x_i \rangle))$$

$$\frac{dL(w)}{dw_j} = \frac{1}{n} \sum_{i=1}^n x_{ij} (\sigma(\langle w, x_i \rangle) - y_i)$$

In matrix form,

$$\nabla L(w) = \frac{1}{n} X^T (\sigma(XW) - Y)$$

**Solution 4b.** The plots are below and the code is attached in following pages.

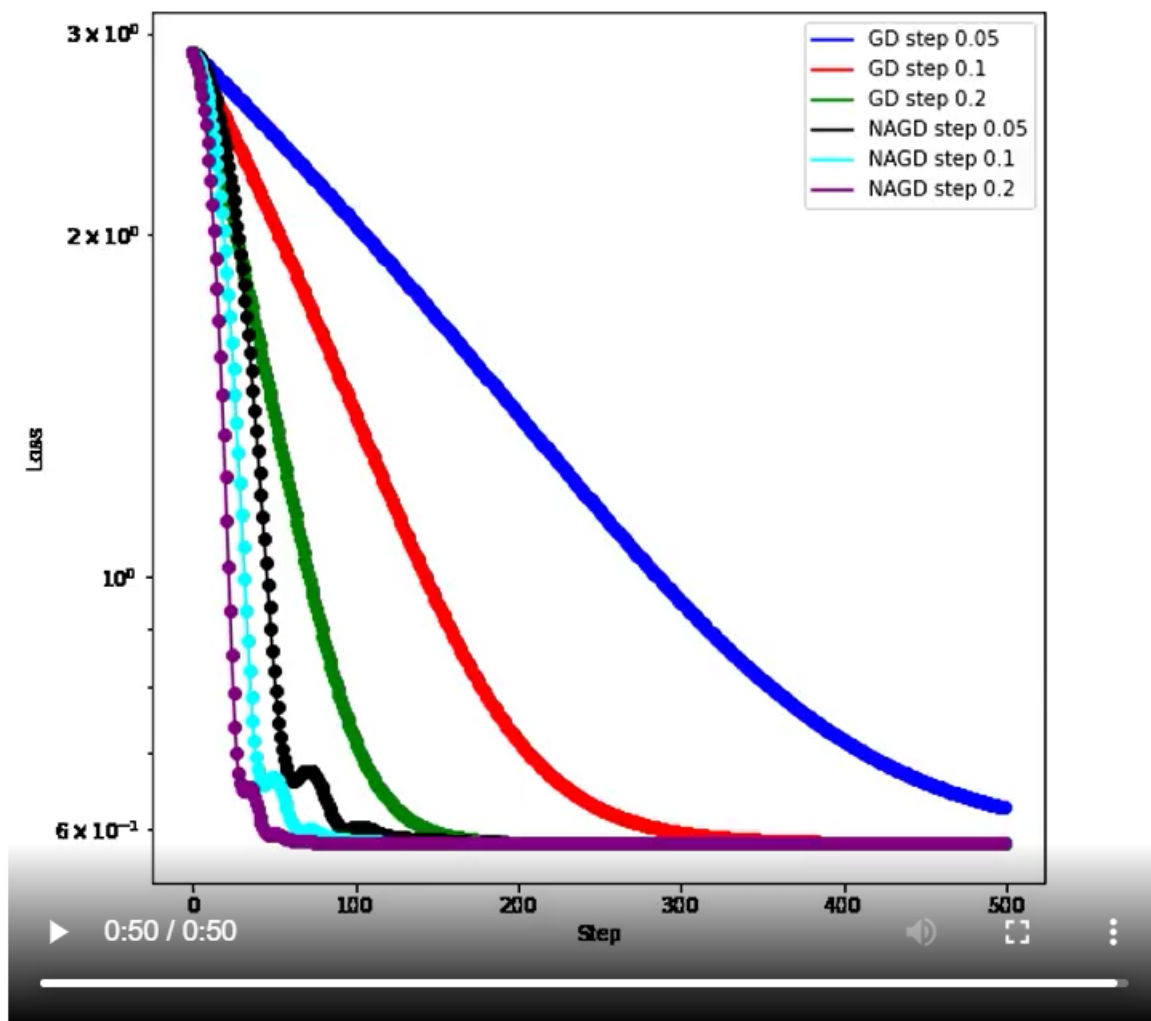


Figure 2: Loss function vs steps

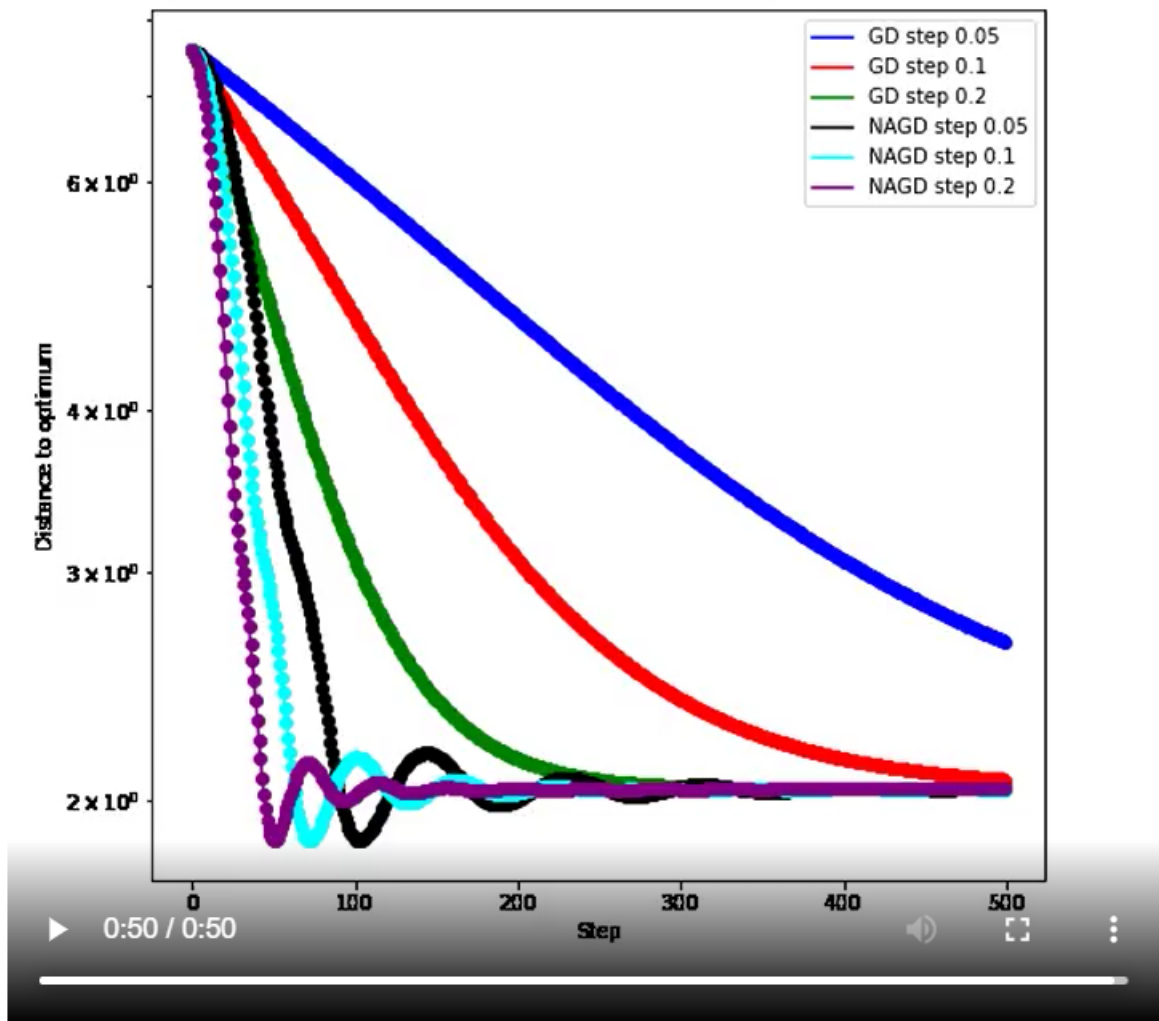


Figure 3: Distance to optimum w vs steps



```
In [1]: import numpy as np
import matplotlib
import matplotlib.pyplot as plt
from matplotlib import cm
from matplotlib.ticker import LinearLocator
from mpl_toolkits.mplot3d import Axes3D
import gzip
from sklearn.preprocessing import OneHotEncoder
from scipy.special import expit
import celluloid
from celluloid import Camera
from matplotlib import animation
from IPython.display import HTML
from matplotlib.lines import Line2D

np.random.seed(2022)
```

```
In [2]: def gradient_descent(xinit,steps,gradient):
        """Run gradient descent.
        Return an array with the rows as the iterates.
        """
        xs = [xinit]
        x = xinit
        for step in steps:
            x = x - step*gradient(x)
            xs.append(x)
        return np.array(xs)

def nagd(winit,gradient,eta=0.1,nsteps=100):
    """Run Nesterov's accelerated gradient descent.
    Return an array with the rows as the iterates.
    """
    ws = [winit]
    u = v = w = winit
    for i in range(nsteps):
        etai = (i+1)*eta/2
        alphai = 2/(i+3)
        w = v - eta*gradient(v)
        u = u - etai*gradient(v)
        v = alphai*u + (1-alphai)*w
        ws.append(w)
    return np.array(ws)
```

```
In [3]: def animated_lplot(Ys,labels=['1','2','3','4','5','6'],ylabel='Function value'):
        """Animated line plot of the Y values.
        Ys is a list where each element is an array of numbers to plot.
        """
        colors = ['blue','red','green','black','cyan','purple','pink']
        fig, ax = plt.subplots(figsize=(6,6))
        camera = Camera(fig)
        T = len(Ys[0])
        plt.yscale('log')
        for t in range(T):
            for j in range(len(Ys)):
                plt.plot(range(t),Ys[j][:t],color=colors[j],marker='o')
            camera.snap()
        handles = []
        for i in range(len(Ys)):
            handles.append(Line2D([0], [0], color=colors[i], label=labels[i]))
        plt.legend(handles = handles, loc = 'upper right')
        plt.xlabel('Step')
        plt.ylabel(ylabel)
        animation = camera.animate(interval=100,blit=False)
        plt.close()
        animation.save('animation.mp4');
        return animation
```

In [9]:

```

def sigmoid(x):
    return 1/(1 + np.exp(-x))

def generate_logistic(n,d):
    """Generate a dataset under the logistic model.
    """
    X = np.random.randn(n,d)
    wstar = np.random.randn(d,1)
    wstar = wstar/np.linalg.norm(wstar)

    preds = sigmoid(np.dot(X,wstar))
    y = np.ones((n,1))
    y[preds > np.random.rand(n,1)] = 0
    return (X,y,wstar)

def logits_loss(a,b):
    return -a*np.log(b) - (1-a)*np.log(1-b)

def logistic_cost(X,Y,w):
    n,d = X.shape
    b = sigmoid(X.dot(w))
    return np.average(logits_loss(Y,b))

def logistic_gradient(X,Y,w):
    n,d = X.shape
    return (1/n)*(X.T.dot(sigmoid(X.dot(w))-Y))

```

In [25]:

```

n,d = 1000, 50
X,Y,wstar = generate_logistic(n,d)
objective = lambda w: logistic_cost(X, Y, w)
gradient = lambda w: logistic_gradient(X, Y, w)

w0 = np.random.randn(d,1)
steps = [0.05, 0.1, 0.2]
wgd = [gradient_descent(w0, [step]*500, gradient) for step in steps]
wnagd = [nagd(w0,gradient,step,500) for step in steps]

HTML(animated_lplot([[objective(w) for w in wgd] for wgd in wgd] + [[objective(w) for w in wn
    ['GD step ' + str(x) for x in steps]+['NAGD step ' + str(x) for x in steps]])

```

Out[25]:

0:49 / 0:50

```
In [26]: distance_gd = [[np.linalg.norm(w - wstar) for w in ws] for ws in wgd]
distance_nagd = [[np.linalg.norm(w - wstar) for w in wnagdi] for wnagdi in wnagd]
animl = animated_lplot(distance_gd+distance_nagd,
                        ['GD step ' + str(x) for x in steps]+'NAGD step ' + str(x) for x in step
                        HTML(animl.to_html5_video()))
```

Out[26]:

0:49 / 0:50

