

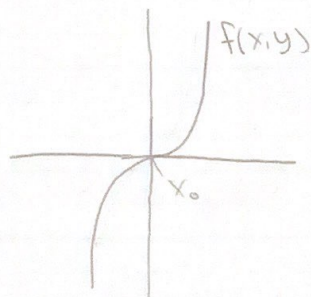
①

$$a.) \quad f(x, y) = x^3 y, \quad x_0 = 0$$

$$\nabla f = \langle 3x^2 y, x^3 \rangle$$

$$\nabla f(x_0) = 0$$

Geographically, we know that despite $\nabla f(x_0) = 0$ for this function, x_0 is not a local minimum or a local maximum



b.) From smoothness upper-bound: f is β -smooth, $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\forall a, h \quad f(a+h) \leq f(a) + \langle \nabla f(a), h \rangle + \frac{\beta}{2} \|\nabla f(a)\|_2^2$$

$$w_{k+1} = w_k - \eta \nabla f(w_k) \quad // \text{ GD}$$

$$\Rightarrow f(w_{k+1}) \leq f(w_k) + \langle \nabla f(w_k), -\eta \nabla f(w_k) \rangle + \frac{\beta}{2} \|\eta \nabla f(w_k)\|_2^2$$

$$\Rightarrow f(w_{k+1}) \leq f(w_k) - \eta \|\nabla f(w_k)\|_2^2 + \frac{\eta^2 \beta}{2} \|\nabla f(w_k)\|_2^2$$

$$\text{Let } \eta = \frac{1}{\beta}$$

$$f(w_{k+1}) \leq f(w_k) - \frac{1}{2\beta} \|\nabla f(w_k)\|_2^2$$

In k GD steps, we have

$$f_* - f(w_0) \leq f(w_k) - f(w_0) \leq \frac{1}{2\beta} \sum_{i=0}^{k-1} \|\nabla f(w_i)\|_2^2$$

when f_* is global optimum

$$\Rightarrow \frac{1}{k} \sum_{i=0}^{k-1} \|\nabla f(w_i)\|_2^2 \leq \frac{2\beta (f(w_0) - f_*)}{k}$$

$$\text{So, for } w_i \text{ we have } \|\nabla f(w_i)\|_2^2 \leq \frac{2\beta (f(w_0) - f_*)}{k}$$

And, to find our point $\|\nabla f(w)\|_2^2 \leq \epsilon$, we need

$$\text{at least } \frac{2\beta (f(x_0) - f_*)}{\epsilon^2} \text{ iterations}$$

② In 1b, we showed that for a β -smooth function f ,

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

For α -convex function that is smooth & $\alpha > 0$, we know

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

$$f(x_k) - f(x^*) \leq \langle \nabla f(x_k), x_k - x^* \rangle - \frac{\alpha}{2} \|x^* - x_k\|_2^2$$

$$= -\frac{1}{2} \|\sqrt{\alpha}(x_k - x^*)\|_2^2 - \frac{1}{\sqrt{\alpha}} \|\nabla f(x_k)\|_2^2 + \frac{1}{2\alpha} \|\nabla f(x_k)\|_2^2$$

$$\leq \frac{1}{2\alpha} \|\nabla f(x_k)\|_2^2$$

Combine inequalities:

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

$$\leq f(x_k) - f(x^*) - \frac{\alpha}{\beta} (f(x_k) - f(x^*))$$

$$= \left(1 - \frac{\alpha}{\beta}\right) (f(x_k) - f(x^*))$$

Iterate this until

$$f(x_k) - f(x^*) \leq \left(1 - \frac{\alpha}{\beta}\right)^k (f(x_0) - f(x^*))$$

$$\Rightarrow \|x_k - x^*\|_2^2 \leq \left(1 - \frac{\alpha}{\beta}\right)^k \|x_0 - x^*\|_2^2$$

⑤ Let x^* be any point that satisfies $\nabla f(x^*) = 0$ & therefore a local minimum

Proof: From convexity

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \forall x, y$$

In particular,

$$f(y) \geq f(x^*) + \langle \nabla f(x^*), y - x^* \rangle \quad \forall y$$

$$\parallel \nabla f(x^*) = 0,$$

$$\Rightarrow f(y) \geq f(x^*) \quad \forall y$$

So, for any y , $f(y) \geq f(x^*)$, proving that

for any point x^* that satisfies $\nabla f(x^*) = 0$, x^* is

a global minimum

④ logit loss function: $l(a, b) = -a \log b - (1-a) \log(1-b)$

$$h_w(x) = \sigma(w \cdot x)$$

$$\sigma(t) = \frac{1}{1+e^{-t}}$$

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

a.)

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

$$\nabla L(w) = \frac{1}{n} \sum_{i=1}^n \nabla_w \left[-y_i \log(\sigma(\langle w, x_i \rangle)) - (1-y_i) \log(1-\sigma(\langle w, x_i \rangle)) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n - \left[\frac{y_i}{\sigma(\langle w, x_i \rangle)} \cdot \frac{\partial}{\partial w} \sigma(\langle w, x_i \rangle) + \frac{1-y_i}{(1-\sigma(\langle w, x_i \rangle))} \cdot \frac{\partial}{\partial w} (1-\sigma(\langle w, x_i \rangle)) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n - \left[\frac{y_i}{\sigma(\langle w, x_i \rangle)} - \frac{1-y_i}{1-\sigma(\langle w, x_i \rangle)} \right] \frac{\partial}{\partial w} \sigma(\langle w, x_i \rangle)$$

$$\parallel \frac{\partial}{\partial w} \sigma(\langle w, x_i \rangle) = -(1+e^{-w \cdot x_i})^{-2} \cdot -x_i = \frac{x_i}{(1+e^{-w \cdot x_i})^2} = \sigma(\langle w, x_i \rangle)(1-\sigma(\langle w, x_i \rangle)) x_i$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n - \left[\frac{y_i - y_i \sigma(\langle w, x_i \rangle) - \sigma(\langle w, x_i \rangle) + y_i \sigma(\langle w, x_i \rangle)}{\sigma(\langle w, x_i \rangle)(1-\sigma(\langle w, x_i \rangle))} \right] \sigma(\langle w, x_i \rangle)(1-\sigma(\langle w, x_i \rangle)) x_i$$

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

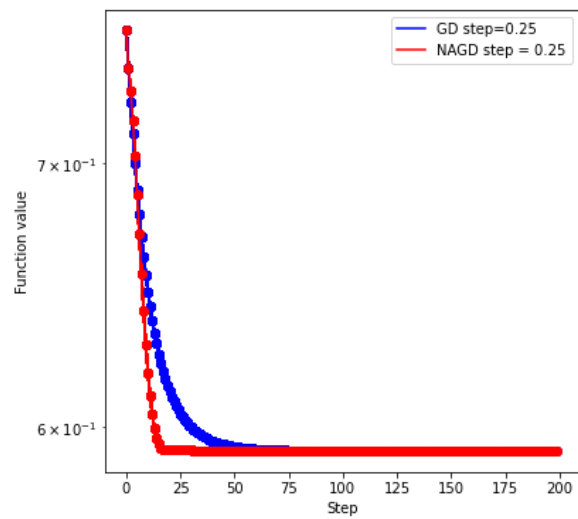
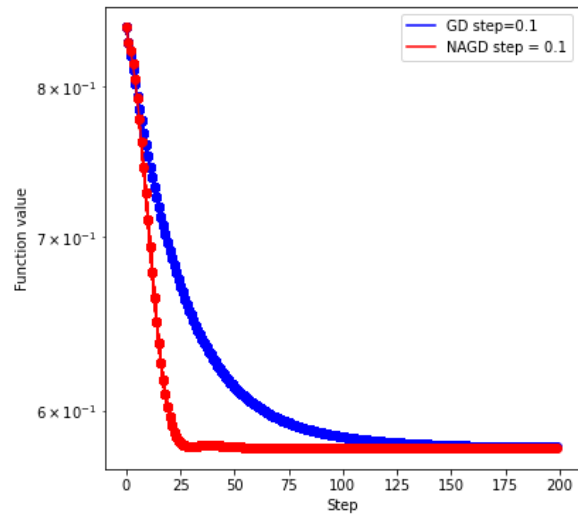
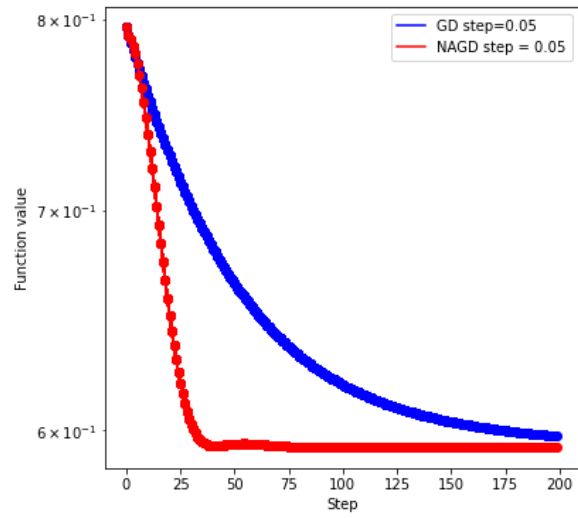
$$\Rightarrow \boxed{\nabla L(w) = \frac{1}{n} \sum_{i=1}^n \left\{ [\sigma(\langle w, x_i \rangle) - y_i] x_i \right\}}$$

b.)

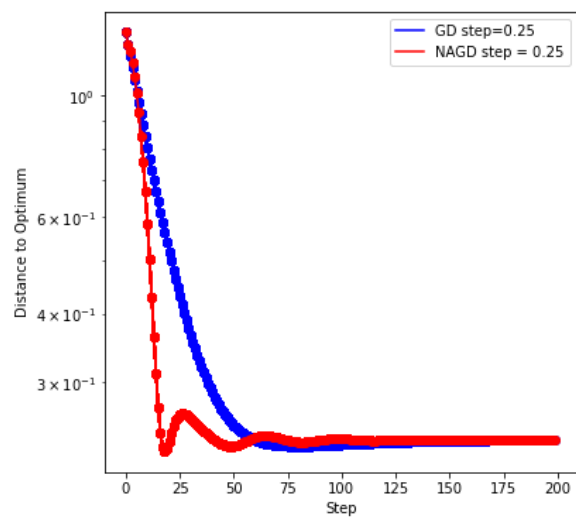
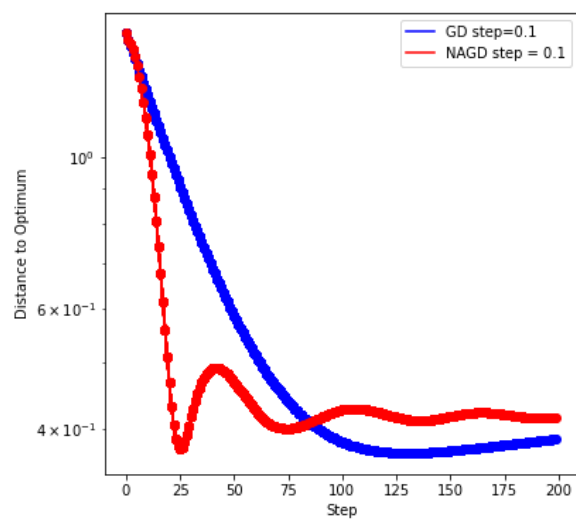
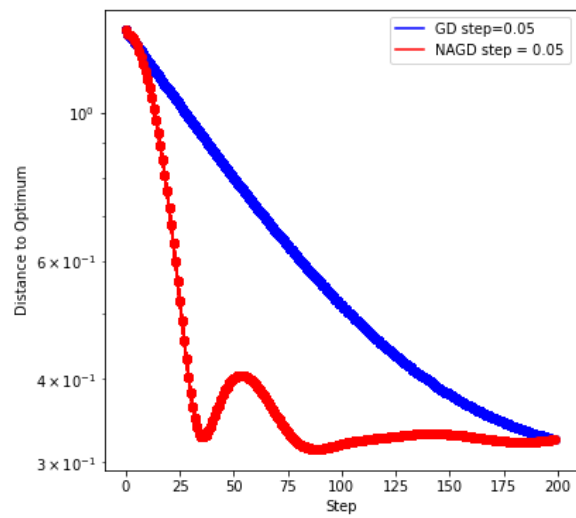
i.)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib.lines import Line2D
4 import time
5
6 def gradient_descent(xinit,steps,gradient):
7     xs = [xinit]
8     x = xinit
9     for step in steps:
10         x = x - step*gradient(x)
11         xs.append(x)
12     return np.array(xs)
13
14 def nagd(winit,gradient,eta,nsteps):
15     ws = [winit]
16     u = v = w = winit
17     for i in range(nsteps):
18         etai = (i+1)*eta/2
19         alpai = 2/(i+3)
20         w = v - etai*gradient(v)
21         u = u - etai*gradient(v)
22         v = alpai*u + (1-alpai)*w
23         ws.append(w)
24     return np.array(ws)
25
26 def plotting(Ys,labels=['1','2','3','4','5','6'],ylabel='Function value'):
27     colors = ['blue','red','green','black','cyan','purple','pink']
28     fig, ax = plt.subplots(figsize=(6,6))
29     T = len(Ys[0])
30     plt.yscale('log')
31     for t in range(T):
32         for j in range(len(Ys)):
33             plt.plot(range(t),Ys[j][:t],color=colors[j],marker='o')
34     handles = []
35     for i in range(len(Ys)):
36         handles.append(Line2D([0],[0],color=colors[i],label=labels[i]))
37     plt.legend(handles = handles, loc = 'upper right')
38     plt.xlabel('Step')
39     plt.ylabel(ylabel)
40
41 def sigmoid(t):
42     return 1/(1+np.exp(-t))
43
44 def y_labeling(p):
45     randint = np.random.uniform(0,1)
46     if randint < p:
47         return 1
48     else:
49         return 0
50
51 def sig_LR(X,Y,w):
52     n,d = X.shape
53     return -(1/n)*(Y.dot(np.log(sigmoid(X.dot(w)))) + (1-Y).dot(np.log(1-sigmoid(X.dot(w)))))
54
55 def sig_LR_gradient(X,Y,w):
56     n,d = X.shape
57     return (1/n)*(X.T.dot(sigmoid(X.dot(w)) - Y))
58
59 n, d = 1000, 20
60 X = np.random.normal(0, 1, (n, d))
61 wstar = np.random.normal(0, 1, d)
62 wstar = wstar/np.linalg.norm(wstar) # unit-norm
63 Y = []
64 for i in range(n): Y.append(y_labeling(sigmoid(X[i].dot(wstar))))
65 Y = np.array(Y)
66 w0 = np.random.normal(0, 1, d)
67 w0 = w0/np.linalg.norm(w0)
68
69 objective = lambda w: sig_LR(X, Y, w)
70 gradient = lambda w: sig_LR_gradient(X, Y, w)
71
72 wgd = gradient_descent(w0, [0.1]*300, gradient)
73 wnagd = nagd(w0,gradient,0.1,300)
74
75 obj_gd = [objective(w) for w in wgd]
76 obj_nagd = [objective(w) for w in wnagd]
77 obj_wstar = [objective(wstar) for w in wgd]
78
79 plotting([obj_gd,obj_nagd],['GD step=0.1','NAGD step = 0.01'])
80
81 distance_gd = [np.linalg.norm(w - wstar) for w in wgd]
82 distance_nagd = [np.linalg.norm(w - wstar) for w in wnagd]
83
84 plotting([distance_gd,distance_nagd],['GD step=0.1','NAGD step = 0.01'],'Distance to Optimum')
85
```

ii.)



iii.)



5.)

```
87  ## Number 5 ##
88
89  def lsr_cost(X,Y,w):
90      n,d = X.shape
91      return (1/n)* np.linalg.norm(X.dot(w) - Y)**2
92
93  def lsr_gradient(X,Y,w):
94      n,d = X.shape
95      return (2/n)*X.T.dot(X.dot(w) - Y)
96
97  n,d = 1000, 1000
98  X = np.random.normal(0, 1, (n, d))
99  wstar = np.random.normal(0, 1, d)
100  noise = np.random.normal(0, 0.1, n)
101  Y = X.dot(wstar) + noise
102  objective = lambda w: lsr_cost(X, Y, w)
103  gradient = lambda w: lsr_gradient(X, Y, w)
104  w0 = np.random.normal(0, 1, d)
105
106  t1 = time.time()
107  ws = gradient_descent(w0, [0.05]*100, gradient)
108  gd = [objective(w) for w in ws]
109  time_gd = time.time() - t1
110
111  t2 = time.time()
112  wnagd = nagd(w0,gradient,0.05,100)
113  nagd = [objective(w) for w in wnagd]
114  time_nagd = time.time() - t2
115
116  t3 = time.time()
117  w_formula = np.linalg.lstsq(X,Y,rcond=None)[0]
118  #w_formula1 = np.linalg.inv(X.T.dot(X)).dot(Y.dot(X))
119  formula = [objective(w_formula) for w in ws]
120  time_formula = time.time() - t3
```

To analyze the time of running GD, NAGD, and the LSR formula, I ran the above code with varying n and d . At small values of n and d , there was not a significant advantage to using GD or NAGD as compared to the least squares regression function. However, as n and d were increased it was apparent that GD and NAGD were advantages over the least squares regression function. Large values of d compared to n showed the greater advantage of using GD and NAGD over the least squares regression function. This is expected when we look at the computational cost of each of these. GD and NAGD are $O(nd)$, while the least squares regression function is $O(nd^2 + d^3)$.