

## Homework 7 for Math 173A - Fall 2025

For coding questions, you must submit your code and the requested answers. **Your code must be present to receive points.**

1. Indicate whether the following functions are strongly convex. If so, find the strong convex parameter. Justify your answer.
  - (a)  $f : \mathbb{R} \rightarrow \mathbb{R}$  with  $f(x) = x$ .
  - (b)  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  with  $f(x) = x_1^2 + x_2^2$ .
  - (c)  $f : \mathbb{R} \rightarrow \mathbb{R}$  with  $f(x) = \log(1 + e^x)$ .
2. Here we will prove the PL-inequality used in class to prove fast convergence for strongly convex functions. Let  $F(x)$  be a strongly convex function with parameter  $c$ , i.e.,  $F$  is  $c$ -strongly convex. Our goal is to show that  $F$  satisfies the  $c$ -PL inequality:

$$F(x) - F(x^*) \leq \frac{1}{2c} \|\nabla F(x)\|^2 \quad \text{for all } x \in \mathbb{R}^d. \quad (1)$$

- (a) Fix  $x \in \mathbb{R}^d$  and define the quadratic function

$$q(y) = F(x) + \nabla F(x)^T(y - x) + \frac{c}{2} \|y - x\|^2.$$

Find the  $y^*$  that minimizes  $q(y)$ .

- (b) Show that  $q(y^*) = F(x) - \frac{1}{2c} \|\nabla F(x)\|^2$
- (c) Use the above to deduce (1).
- (d) Explain the proof technique in your own words to demonstrate understanding of what we did.
3. Consider Q3 from HW6, we plan to analyze gradient descent on this function, and want to understand in theory how fast GD will converge. In the end, you will be asked to compare it numerically with the accelerated GDs and with the conjugate gradient method.
  - (a) Check whether  $f$  is  $L$ -smooth and  $\alpha$ -strongly convex for some  $L, \alpha > 0$ . Find the corresponding  $L$  if  $f$  is  $L$ -smooth and the  $\alpha$  if  $f$  is strongly convex.

- (b) Give a range of step sizes, i.e., an interval in this case, so the method will converge using the theorem we learned for strongly convex and smooth functions.
- (c) Continuing (b), what can we say about  $f(x^{(t)}) - f(x^*)$  and its rate of decay according to the theorem we learned? Does it depend on the choice of the stepsize? Which stepsize gives you the best rate? If you want to achieve  $f(x^{(t)}) - f(x^*) \leq 10^{-5}$  for  $n = 100$  and  $f(x^{(0)}) - f(x^*) = 1$ , how many iterations do you need using the best rate you found?
- (d) **Coding Question.** Program and run the conjugate gradient descent method for  $f$  starting at  $x^{(0)} = 0$  for  $n = 20$  and  $n = 100$ . Run the method for  $n$  iterations. Plot the  $f(x^{(t)}) - f(x^*)$  for the conjugate gradient descent method and the three methods you have explored in Q3 from HW6 in the same figure. In a different figure, plot  $\|x^{(t)} - x^*\|$  for the four methods. If you encounter a number smaller than  $10^{-16}$ , set it to be  $10^{-16}$ . In both plots, make the logarithmic scale for the vertical axis. Comment on the plots.
4. Consider the following set in  $\mathbb{R}^n$  for an integer  $s > 0$ :

$$B = \{x \in \mathbb{R}^n \mid x_i \geq 0, \text{ for } i = 1, \dots, n \text{ and } x \text{ has at most } s \text{ nonzeros.}\}.$$

- (a) Find an expression for the orthogonal projection of a point  $x \in \mathbb{R}^n$  onto  $B$  (No need for justification).
- (b) For the function

$$f(x) = \frac{1}{2} \|Ax - b\|^2, \quad (2)$$

Write a projected gradient descent algorithm to solve

$$\min_{x \in \Omega} f(x) \quad (3)$$

for  $\Omega = B$ , with  $B$  from part (a). You need to specify the gradient formula and the projection formula. You do not need to specify the step size for this problem.

- (c) **Coding Question:** Consider the optimization problem (3), where  $A \in \mathbb{R}^{20 \times 50}$  and  $b \in \mathbb{R}^{20}$  are from the dataset `HW7Ab.csv`. The file `HW7Ab.csv` contains the data  $A$  and  $b$ . The first 50 columns form the matrix  $A$  and the last column is the vector  $b$ . The vector  $b$  is generated by setting  $b = Ax^*$  for a vector  $x^* \in \mathbb{R}^{50}$  that has 2 positive nonzeros. Note the linear system  $Ax = b$  is underdetermined and has a lot of solutions. Program the projected gradient method for solving (3) with  $s = 2$  to find the  $x^*$ . You can experiment with the stepsize and initialization to make sure  $f(x^{(t)})$  converges to 0. You need to submit the code, the plot of  $f(x^{(t)}) - f(x^*) = f(x^{(t)})$  vs iteration number, and the indices and values that has absolute value larger than  $10^{-6}$  of  $x^*$  you found.

5. **Coding Question:** We will implement the SVM algorithm with gradient descent to classify two Gaussians in 2D. The dataset is given in `HW7SVM.csv`.

- (a) In `HW7SVM.csv`, the first 100 rows are the data for cluster 1:  $(x_i, y_i) \in \mathbb{R}^2 \times \mathbb{R}$ ,  $i = 1, \dots, 100$ , with  $y_i = 1$  always. The next 100 rows are the data for cluster 2:  $(x_i, y_i) \in \mathbb{R}^2 \times \mathbb{R}$ ,  $i = 101, \dots, 200$ , with  $y_i = -1$  always. Create and turn in a scatter plot of the feature vectors, i.e., the  $x_i$ s, colored by the label, i.e.,  $y_i$ s (blue for 1 and red for -1).
- (b) Create a function for the gradient of the loss

$$L(w) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \max(0, 1 - y_i \langle x_i, w \rangle)$$

$$\nabla L(w) = w + \sum_{i=1}^n -y_i x_i \cdot \mathbf{1}_{1-y_i \langle x_i, w \rangle > 0},$$

where  $\mathbf{1}_{1-y_i \langle x_i, w \rangle > 0} = \begin{cases} 1 & \text{if } 1 - y_i \langle x_i, w \rangle > 0 \\ 0 & \text{else} \end{cases}$ . Also, here  $n = 200$ . To

compute the gradient, you'll have to compute an indicator of whether  $1 - y_i \langle x_i, w \rangle$  is positive or negative at every point, and sum up the contribution of this term for all points where it's positive.

- (c) Setting the step size  $\mu = 10^{-4}$  and starting at  $w^{(0)} = (-1, 1)$ , run 1000 iterations of gradient descent. You will create two plots.
  - i. Plot the classification error (averaged over all the points) as a function of the iterations. The classification of  $x_i$  is determined by  $\text{sign}(\langle x_i, w \rangle)$ .
  - ii. Plot the margin  $\frac{2}{\|w\|}$  as a function of the iterations. This shows how much of a gap you have between the classes you've learned.
- (d) Create another scatter plot of your data, but this time color the points by the function  $f(x_i) = 1 - y_i \cdot \langle x_i, w \rangle$ . The numbers closest to 0 (positive numbers or largest negative numbers) will show you which points were “most important” in determining the classification.

**Note:** Here we're only defining a subspace classifier (i.e. the classifier goes through the origin). This is fine for our problem as the gaussians are on opposite sides of the origin. If you want to create an intercept term, simply append a vector of all 1's as a column of your data, and now your weight vector will be of dimension 3 instead of 2. This is done the same way as when running least squares regression.

1) Indicate whether the following functions are strongly convex.  
If so, find the strong convex parameter. Justify your answer.

Def: A function  $F$  is strongly convex if  $\exists c > 0$  s.t.  $\forall x, y \in \mathbb{R}^d$ ,

$$F(y) \geq F(x) + \nabla F(x)^T (y - x) + \frac{1}{2} c \|y - x\|^2$$

Wmna: A  $C^2$  function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is  $c$ -strongly convex iff

$$\forall x, \nabla^2 f(x) - c \cdot I \leq 0$$

$$\Leftrightarrow \forall x, v, v^T \nabla^2 f(x) v \geq c \|v\|^2$$

$$\Leftrightarrow \forall x, \lambda_{\min}(\nabla^2 f(x)) \geq c \quad \text{for } c > 0$$

(a)  $f: \mathbb{R} \rightarrow \mathbb{R}$  with  $f(x) = x$

Check:  $\forall x, f''(x) \geq c$ , where  $c > 0$

$f''(x) = 0$ , so  $\exists c > 0$  such that  $f''(x) \geq c$

So  $f$  is not strongly convex.

(b)  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  with  $f(x) = x_1^2 + x_2^2$

Check:  $\forall x, \lambda_{\min}(\nabla^2 f(x)) \geq c$

$$\nabla f(x) = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix} \quad \nabla^2 f(x) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\lambda_{\min}(\nabla^2 f(x)) = 2 \Rightarrow c = 2$$

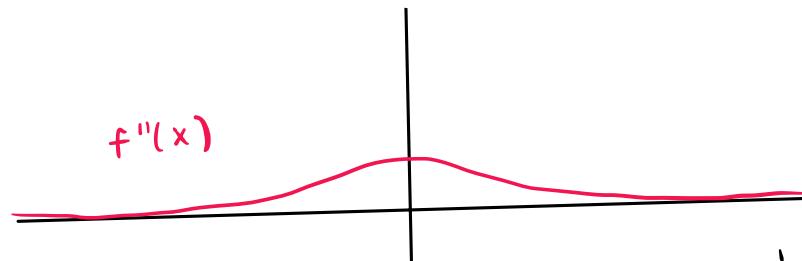
Yes,  $f$  is strongly convex with  $c = 2$ .

(c)  $f: \mathbb{R} \rightarrow \mathbb{R}$  with  $f(x) = \log(1 + e^x)$

$$f'(x) = \frac{e^x}{1+e^x}, \quad f''(x) = \frac{e^x(1+e^x) - e^x \cdot e^x}{(1+e^x)^2} = \frac{e^x}{(1+e^x)^2} > 0$$

Want  $\frac{e^x}{(1+e^x)^2} \geq c$ , where  $c > 0$ , for all  $x$ .

Graph of  $f''(x)$ :



As  $x \rightarrow \pm\infty$ ,  $f''(x) \rightarrow 0$  but never actually gets to it (since  $e^x > 0$  for any  $x$ ). Therefore we cannot find some  $c > 0$  such that  $f''(x) \geq c$  for all values of  $x$ .

So  $f$  is not strongly convex.

2) Here we will prove the PL-inequality used in class to prove fast convergence for strongly convex functions. Let  $F(x)$  be a strongly convex function with parameter  $c$ , i.e.,  $F$  is  $c$ -strongly convex. Our goal is to show that  $F$  satisfies the  $c$ -PL inequality:

$$F(x) - F(x^*) \leq \frac{1}{2c} \|\nabla F(x)\|^2 \quad \forall x \in \mathbb{R}^d$$

(a) Fix  $x \in \mathbb{R}^d$  and define the quadratic function

$$q(y) = F(x) + \nabla F(x)^T (y - x) + \frac{c}{2} \|x - y\|^2$$

Find the  $y^*$  that minimizes  $q(y)$ .

$$\begin{aligned} \nabla q(y) &= \nabla \left( F(x) + \nabla F(x)^T (y - x) + \frac{c}{2} \|x - y\|^2 \right) \\ &= \nabla \left( \nabla F(x)^T y - \cancel{\nabla F(x)^T x} + \frac{c}{2} \|x - y\|^2 \right) \\ &= \nabla F(x) + \nabla \left( \frac{c}{2} \|x - y\|^2 \right) \end{aligned}$$

$$= \nabla F(x) + c(y-x)$$

$$\nabla F(x) + c(y-x) = 0 \Rightarrow \nabla F(x) + cy - cx = 0$$

$$\Rightarrow cy = cx - \nabla F(x) \Rightarrow y^* = x - \frac{1}{c} \nabla F(x)$$

(b) Show that  $q(y^*) = F(x) - \frac{1}{2c} \|\nabla F(x)\|^2$

$$q(y^*) = q\left(x - \frac{1}{c} \nabla F(x)\right)$$

$$= F(x) + \nabla F(x)^T \left(x - \frac{1}{c} \nabla F(x)\right) + \frac{c}{2} \|x - \frac{1}{c} \nabla F(x)\|^2$$

$$= F(x) - \frac{1}{c} \nabla F(x)^T \nabla F(x) + \frac{c}{2} \left\| -\frac{1}{c} \nabla F(x) \right\|^2$$

$$= F(x) - \frac{1}{c} \|\nabla F(x)\|^2 + \left(\frac{c}{2} \cdot \frac{1}{c^2}\right) \|\nabla F(x)\|^2$$

$$= F(x) + \left(-\frac{1}{c} + \frac{1}{2c}\right) \|\nabla F(x)\|^2 = F(x) - \frac{1}{2c} \|\nabla F(x)\|^2$$

(c) Use the above to deduce (1).

$$(1) F(x) - F(x^*) \leq \frac{1}{2c} \|\nabla F(x)\|^2$$

We know that  $F$  is a strongly convex function, thus, by a previous theorem,

$$F(y) \geq F(x) + \nabla F(x)^T (y-x) + \frac{1}{2} c \|y-x\|^2 = q(y)$$

$$\Rightarrow F(y) \geq q(y).$$

We know that  $x^*$  minimizes  $F$ , and above implies

$$F(x^*) \geq q(x^*)$$

We know that  $y^*$  minimizes  $q$ , so

$$F(x^*) \geq q(x^*) \geq q(y^*)$$

$$\Rightarrow F(x^*) \geq q(y^*)$$

$$\Rightarrow F(x^*) \geq F(x) - \frac{1}{2c} \|\nabla F(x)\|^2$$

Rearrange to get

$$F(x) - F(x^*) \leq \frac{1}{2c} \|\nabla F(x)\|^2.$$

(d) Explain the proof technique in your own words to demonstrate understanding of what we did.

In (a) & (b), we found the minimizer of the strong convexity lower bound  $q(y)$  for  $F(y)$ , and then found the minimum value of the lower bound  $q(y^*)$ . Then, we took the minimum value of  $F$ ,  $F(x^*)$ , and used the simple logic of minimizers to deduce that  $F(x^*) \geq q(y^*)$ . Knowing the formula for  $q(y^*)$  from (b), we proved the inequality.

3) Consider Q3 from HW6, we plan to analyze GD on this function, and want to understand in theory how fast GD will converge. In the end, you will be asked to compare it numerically with the accelerated GDs and with the conjugate gradient method.

$A \in \mathbb{R}^{n \times n}$  diagonal matrix with diagonal entries  $A_{ii} = i$

$b \in \mathbb{R}^n$  vector with all 1 entries. Define

$$f(x) = \frac{1}{2} x^T A x - b^T x.$$

(a) Check whether  $f$  is  $L$ -smooth and  $\alpha$ -strongly convex for some  $L, \alpha > 0$ . Find the corresponding  $L$  if  $f$  is  $L$ -smooth and the  $\alpha$  if  $f$  is strongly convex.

L-smooth:  $f$  is  $C^2$ , so we use  $\|\nabla^2 f(x)\|_{op} \leq L \forall x$

$$\nabla f(x) = Ax - b \quad \nabla^2 f(x) = A$$

$$\|A\|_{op} = \lambda_{\max}(A) = n$$

Yes,  $f$  is  $L$ -smooth with  $L = n$ .

Strongly convex:  $\forall x, \lambda_{\min}(\nabla^2 f(x)) \geq \alpha > 0$

$$\lambda_{\min}(A) = 1 \Rightarrow \alpha = 1$$

Yes,  $f$  is  $\alpha$ -strongly convex with  $\alpha = 1$ .

(b) give a range of step sizes, i.e., an interval in this case, so the method will converge using the theorem we learned for strongly convex and smooth functions.

Theorem: If  $F$  is  $\alpha$ -strongly convex and  $L$ -smooth, then

$$F(x^{(t)}) - F(x^*) \leq (1 - \mu\alpha)^t (F(x^{(0)}) - F(x^*))$$

for any  $\mu \in (0, \frac{1}{L}]$

$$\Rightarrow \mu \in (0, \frac{1}{n}]$$

(c) Continuing (b), what can we say about  $f(x^{(t)}) - f(x^*)$  and its rate of decay according to the theorem we learned? Does it depend on the choice of the stepsize? Which stepsize gives you the best rate? If you want to achieve  $f(x^{(t)}) - f(x^*) \leq 10^{-5}$  for  $n=100$  and  $f(x^{(0)}) - f(x^*) = 1$ , how many iterations do you need using the best rate you found?

From theorem, with our  $\alpha$  and  $L$ :

$$f(x^{(t)}) - f(x^*) \leq (1 - \mu)^t (f(x^{(0)}) - f(x^*)) , \mu \in (0, \frac{1}{n}]$$

The rate of decay definitely depends on  $\mu$ , since at each iteration the decay increases by a factor of  $(1-\mu)$ . The "best" rate of decay would mean faster, so we'd want  $(1-\mu)$  to be larger. Since  $\mu \in [0, \frac{1}{n}]$ , we want  $\mu$  to be the "smallest" in that range so when we subtract from 1 it gives the largest possible value. So  $\mu = \frac{1}{n}$  would give the "best" rate.

$$\text{Want: } f(x^{(t)}) - f(x^*) \leq 10^{-5} \text{ with } n=100, f(x^{(0)}) - f(x^*) = 1$$

$$\Rightarrow f(x^{(t)}) - f(x^*) \leq \left(1 - \frac{1}{100}\right)^t \leq 10^{-5}$$

$$\Rightarrow (0.99)^t \leq 10^{-5}$$

$$\Rightarrow \ln((0.99)^t) \leq \ln(10^{-5})$$

$$\Rightarrow t \geq \frac{\ln(10^{-5})}{\ln(0.99)}$$

$$\Rightarrow t \geq \frac{-11.51}{-0.01}$$

$$\Rightarrow t \geq 1145.53 \Rightarrow t = 1146 \text{ iterations}$$

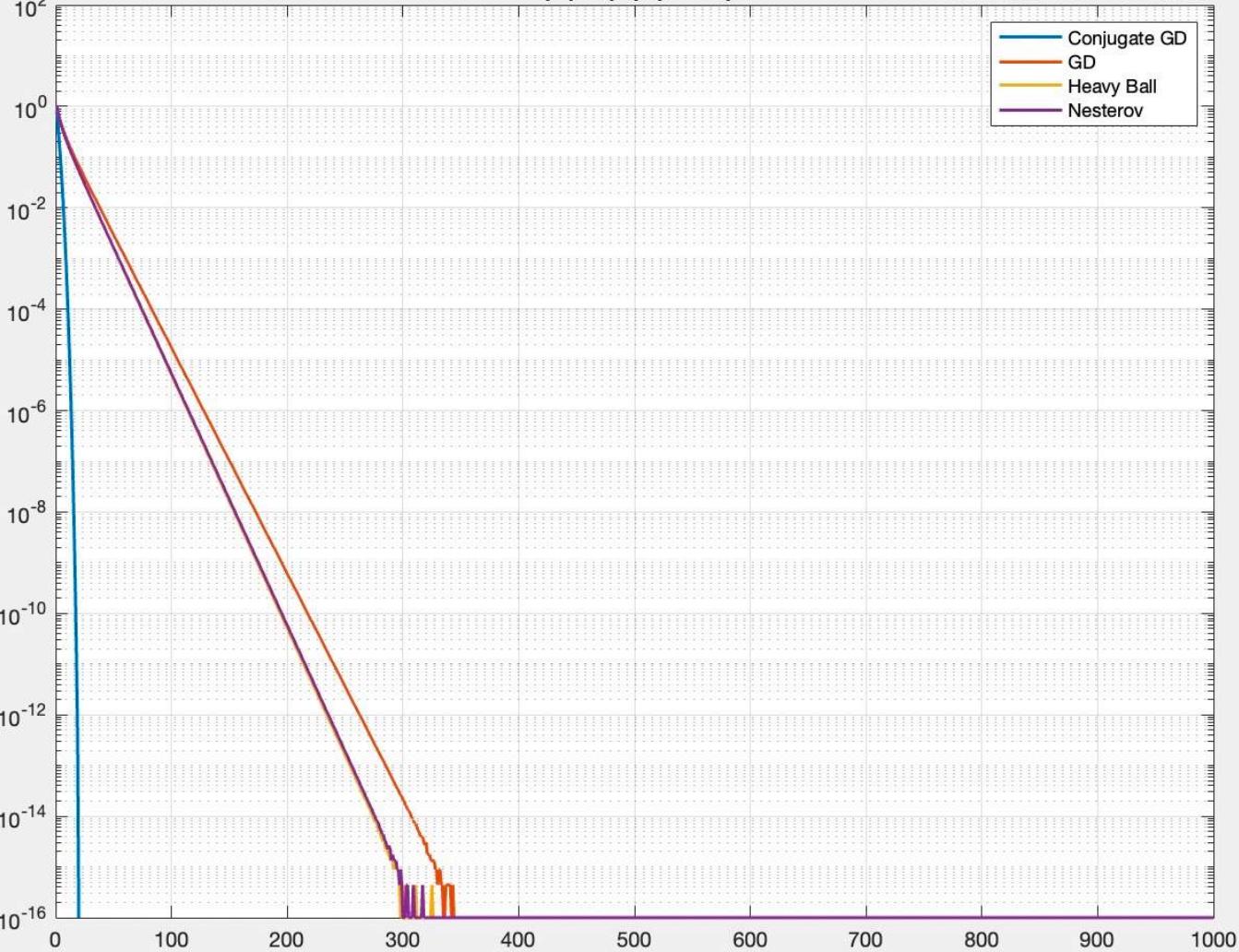
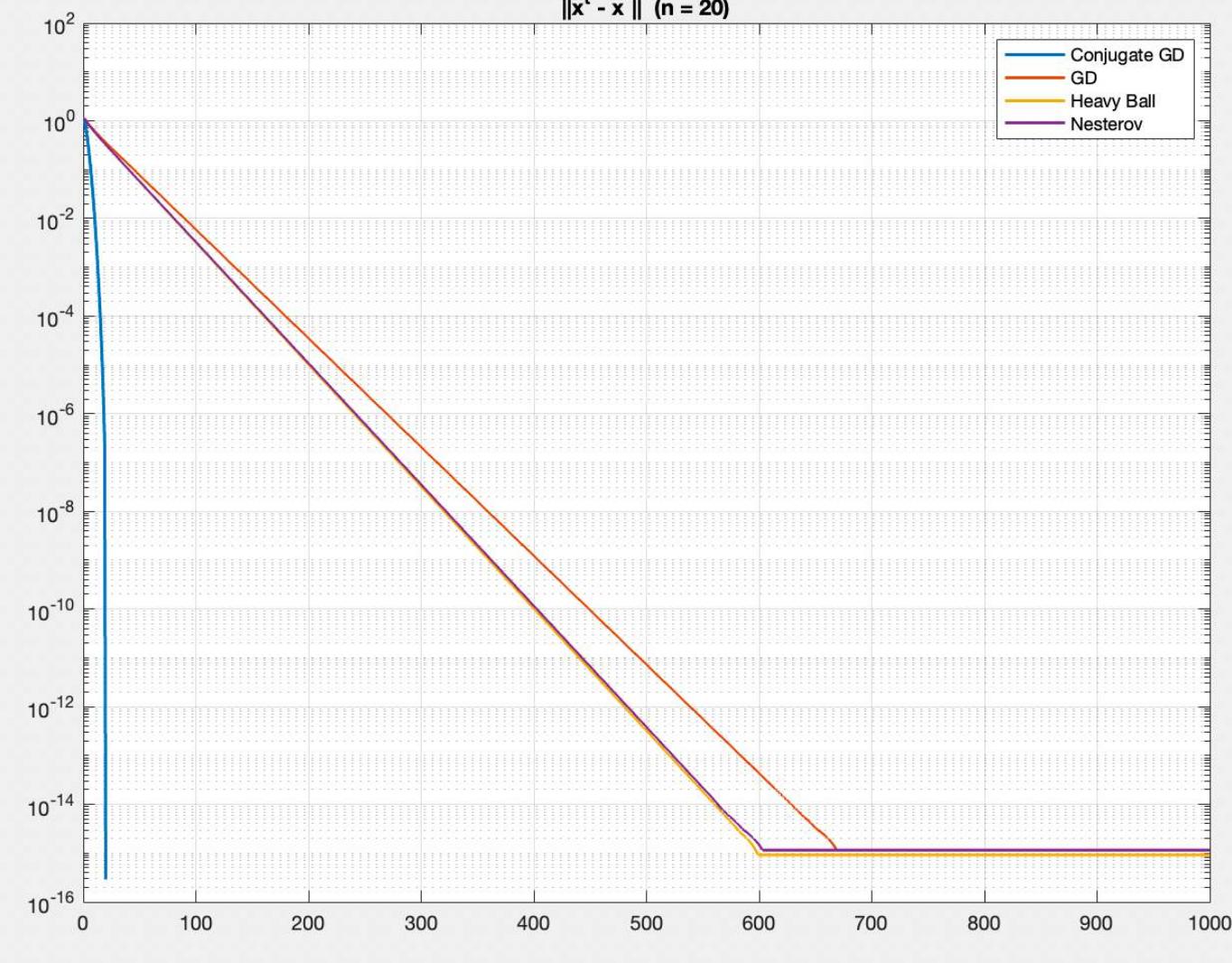
(d) Coding Question. Program and run the conjugate gradient descent method for  $f$  starting at  $x^{(0)} = 0$  for  $n = 20$  and  $n = 100$ . Run the method for  $n$  iterations. Plot the  $f(x^{(t)}) - f(x^*)$  for the conjugate gradient descent method and the three methods you have explored in Q3 from HW6 in the same figure. In a different figure, plot  $\|x^{(t)} - x^*\|$  for the four methods. If you encounter a number smaller than  $10^{-10}$ , set it to be  $10^{-10}$ . In both plots, make the logarithmic scale for the vertical axis. Comment on the plots.

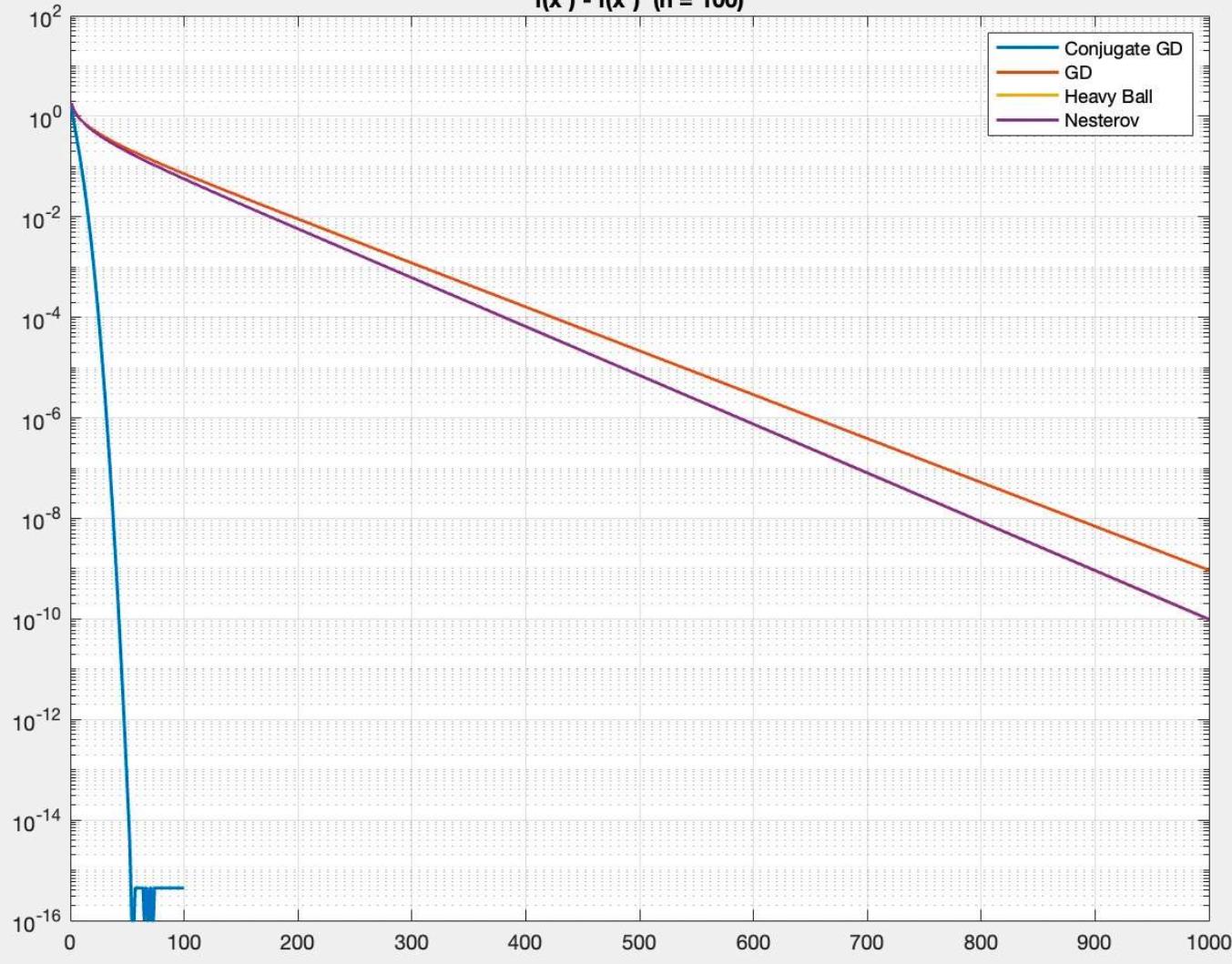
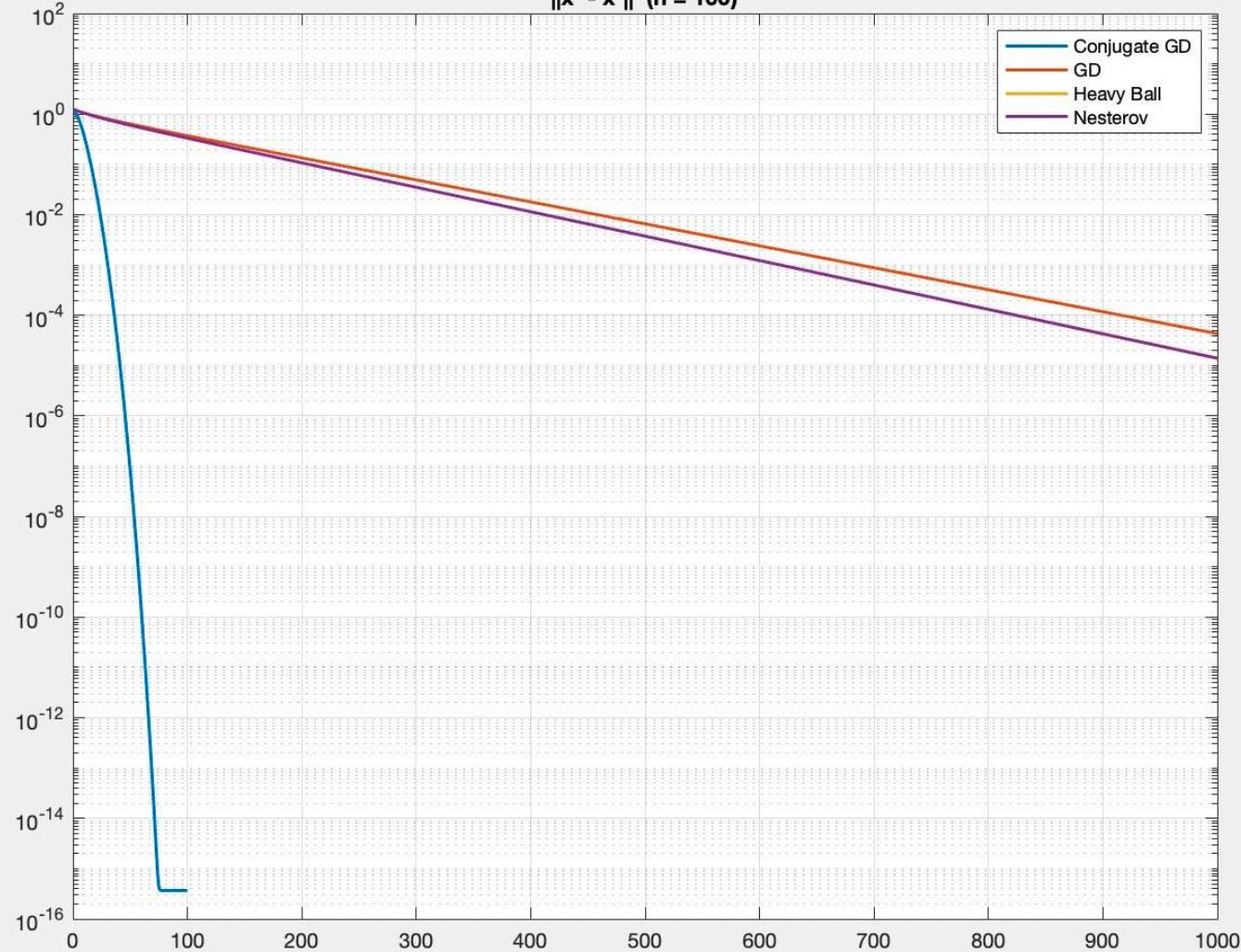
```

function [f_error, x_error] = conjugate_GD(n)
    A = diag(1:n);
    b = ones(n, 1);
    x0 = zeros(n, 1);
    x_star = A\b;
    f_star = 0.5 * x_star' * A * x_star - b' * x_star;
    f_error = zeros(n,1);
    x_error = zeros(n, 1);
    x = x0;
    r = b - A*x;
    p = r;
    r_old = r' * r;
    for k = 1:n
        Ap = A * p;
        alpha = r_old /(p'*Ap);
        x = x + alpha * p;
        r = r - alpha * Ap;
        r_new = r' * r;
        f_val = 0.5 * x' * A * x - b' * x;
        f_error(k) = f_val - f_star;
        x_error(k) = norm(x - x_star);
        beta = r_new / r_old;
        p = r + beta*p;
        r_old = r_new;
    end
end

1   for n = [20, 100]
2       A = diag(1:n);
3       b = ones(n,1);
4       x_star = A\b;
5
6       [f(CG, x(CG)] = conjugate_GD(n);
7       [f(GD, x(GD)] = gd_hw6(n);
8       | [f(HB, x(HB)] = gd_HB(n);
9       | [f(NA, x(NA)] = gd_NA_2(n);
10
11      f(CG(f(CG < 1e-16) = 1e-16;
12      f(GD(f(GD < 1e-16) = 1e-16;
13      f(HB(f(HB < 1e-16) = 1e-16;
14      f(NA(f(NA < 1e-16) = 1e-16;
15
16      x(CG(x(CG < 1e-16) = 1e-16;
17      x(GD(x(GD < 1e-16) = 1e-16;
18      x(HB(x(HB < 1e-16) = 1e-16;
19      x(NA(x(NA < 1e-16) = 1e-16;
20
21      figure;
22      semilogy(f(CG, 'LineWidth', 1.5); hold on;
23      semilogy(f(GD, 'LineWidth', 1.5);
24      semilogy(f(HB, 'LineWidth', 1.5);
25      semilogy(f(NA, 'LineWidth', 1.5);
26      title(sprintf("f(x^t) - f(x^*) (n = %d)", n));
27      legend("Conjugate GD", "GD", "Heavy Ball", "Nesterov");
28      grid on;
29
30      figure;
31      semilogy(x(CG, 'LineWidth', 1.5); hold on;
32      semilogy(x(GD, 'LineWidth', 1.5);
33      semilogy(x(HB, 'LineWidth', 1.5);
34      semilogy(x(NA, 'LineWidth', 1.5);
35      title(sprintf("||x^t - x^*|| (n = %d)", n));
36      legend("Conjugate GD", "GD", "Heavy Ball", "Nesterov");
37      grid on;
38
39  end

```

$f(x^t) - f(x^*)$  ( $n = 20$ ) $\|x^t - x^*\|$  ( $n = 20$ )

$f(x^t) - f(x^*)$  ( $n = 100$ ) $\|x^t - x^*\|$  ( $n = 100$ )

n=20: We can see from the plot that conjugate GD has a very fast convergence compared to normal GD, HB, and NA, which all behave very similarly. In the  $f(x^{(t)}) - f(x^*)$  plot, conjugate GD appears to converge to  $10^{-16}$  in exactly n iterations, while GD, HB, and NA all take  $\sim 300+$  iterations to converge. In this plot, GD, NA, and HB oscillate after reaching convergence, which I think is because I did not tell any of the methods to stop running once they converged. We can still see that regular GD has the slowest convergence, HB/NA have similar rates that are slightly faster than GD, and conjugate GD converges much, much faster than the rest. The same can be said for the  $\|x^{(t)} - x^*\|$  plot about the rates of convergence, however, this time GD, HB, NA converge in about 600+ iterations, so the difference between conjugate GD is even greater. There is no oscillation in this plot, so perhaps the oscillation in  $f(x^{(t)}) - f(x^*)$  had to do with the stepsize choices being optimal for finding  $x^*$ .

n=100: Again, in both  $f(x^{(t)}) - f(x^*)$  and  $\|x^{(t)} - x^*\|$ , we see conjugate GD converge much, much faster than GD, HB, and NA. Conjugate GD actually appears to converge in less than n iterations, oscillating to the  $10^{-16}$  bound due to fast convergence. Compared to n=20, the GD, HB, and NA appear to converge much, much slower, still with similar rates to one another. So, conjugate GD runs "better" with larger n, and GD, HB, NA run "worse" with larger n, with their step sizes being inversely proportional to n.

4) Consider the following set in  $\mathbb{R}^n$  for an integer  $s > 0$ :

$$B = \{x \in \mathbb{R}^n \mid x_i \geq 0, \text{ for } i=1, \dots, n \text{ and } x \text{ has at most } s \text{ nonzeros}\}.$$

(a) Find an expression for the orthogonal projection of a point  $x \in \mathbb{R}^n$  onto  $B$  (NO NEED TO JUSTIFY)

$$(\Pi_B(x \in \mathbb{R}^n))_i = \begin{cases} \max(x_i, 0), & x_i \in \{s \text{ largest entries of } x\} \\ 0, & x_i \notin \{s \text{ largest entries of } x\} \end{cases}$$

(b) For the function  $f(x) = \frac{1}{2} \|Ax - b\|^2$ , write a projected gradient descent alg to solve  $\min_{x \in \mathbb{R}^n} f(x)$  for  $\mathcal{N} = B$ . You need to specify the gradient formula and the projection formula. You do not need to specify step size.

alg proj-GD

Initialize  $x^{(0)} \in B$

(Initialize stepsize - not specified in problem)

Initialize  $s > 0$

for  $i = 0, 1, \dots, d$  do

$$\text{grad} = A^T(Ax^{(i)} - b)$$

(update stepsize - not specified in this problem)

$$y^{(i+1)} = x^{(i)} - (\text{stepsize}) * \text{grad}$$

for  $j = 1, \dots, n$  do

$$\text{if } y^{(i)}[j] < 0 : y^{(i)} = 0 \text{ end}$$

end

$$x^{(i+1)} = \text{zeros}(n)$$

$$\text{count} = 0$$

Initialize empty list L

```

for j = 1, ..., n do
    if  $y^{(i)}[j] > 0$  : count ++
    end
    add  $y^{(i)}[j]$  to L
end
if count ≤ s :  $x^{(i+1)} = y^{(i)}$  end
else :
    Sort L by descending value
    For j = 1, ..., s do
        value = L[j-1]
         $x^{(i+1)}[j-1] = value$ 
    end
end
end

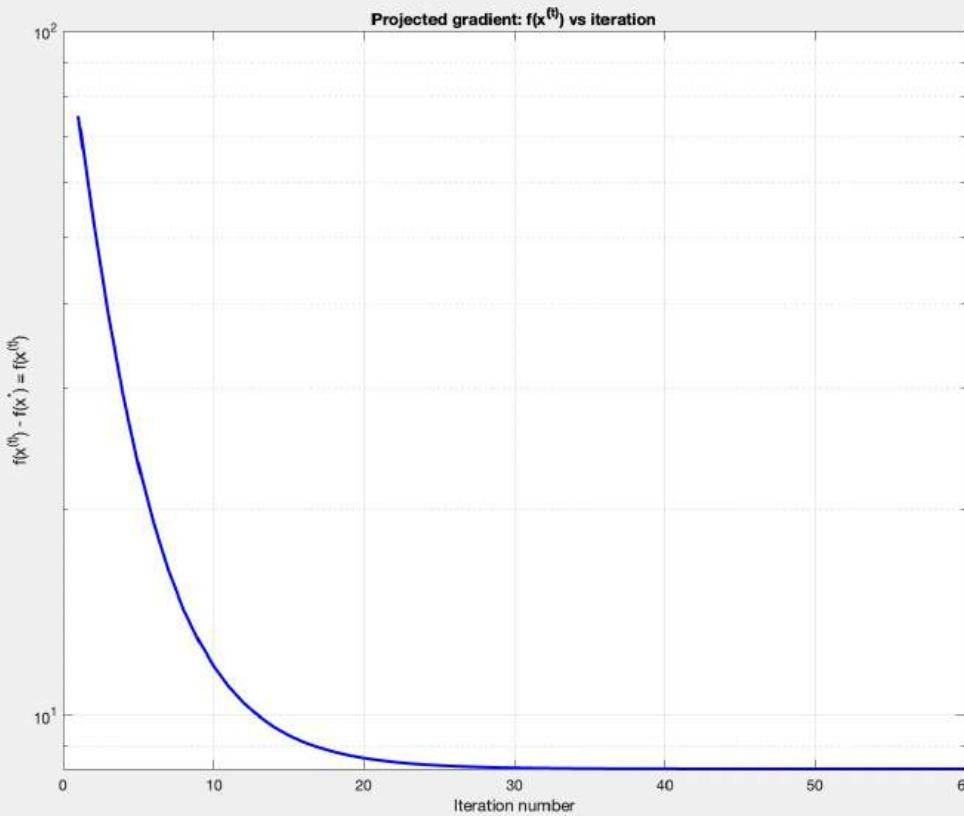
```

(c) Coding Question. Consider problem (3), where  $A \in \mathbb{R}^{20 \times 50}$  and  $b \in \mathbb{R}^{20}$  are from dataset HW7Ab.csv. The csv file contains the data A and b. First 50 columns form A, last one is b. The vector b is generated by setting  $b = Ax^*$  for a vector  $x^* \in \mathbb{R}^{50}$  that has 2 positive nonzeros. Note  $Ax = b$  is underdetermined and has a lot of solutions. Program the projected gradient method for solving (3) with  $\delta = 2$  to find the  $x^*$ . You can experiment with the stepsize and initialization to make sure  $f(x^{(t)})$  converges to 0. You need to submit the code, the plot  $f(x^{(t)}) - f(x^*) = f(x^{(t)})$  vs iteration number, and the indices and values that has abs val larger than  $10^{-6}$  of  $x^*$  you found.

```

1 function x = HW7_q4()
2     data = readmatrix('HW7Ab.csv');
3     A = data(:, 1:50);
4     b = data(:, 51);
5     [~, n] = size(A);
6     s = 2;
7     x0 = zeros(n, 1);
8     x = x0;
9     max_iter = 60;
10    L = norm(A'*A, 2);
11    alpha = 1/L;
12    f_vals = zeros(max_iter, 1);
13    x_error = zeros(max_iter, 1);
14    for i = 1:max_iter
15        grad = A' * (A*x-b);
16        f_vals(i) = 0.5 * norm(A*x-b)^2;
17        y = x - alpha * grad;
18        y = max(y, 0);
19        [sorted_vals, sorted_idx] = sort(y, 'descend');
20        x_new = zeros(n, 1);
21        keep = sorted_idx(1:min(s, length(sorted_idx)));
22        x_new(keep) = y(keep);
23        x_error(i) = norm(x_new - x);
24        x = x_new;
25    end
26    above_threshold = abs(x) > 1e-6;
27    if sum(above_threshold) > 0
28        fprintf('Index      Value\n');
29        for i = 1:n
30            if above_threshold(i)
31                fprintf('%3d      %12.8f\n', i, x(i));
32            end
33        end
34    end
35    figure;
36    semilogy(1:length(f_vals), f_vals, 'b-', 'LineWidth', 2);
37    xlabel('Iteration number');
38    ylabel('f(x^{(t)}) - f(x^{(*)}) = f(x^{(t)})');
39    title('Projected gradient: f(x^{(t)}) vs iteration');
40    grid on;
41 end

```



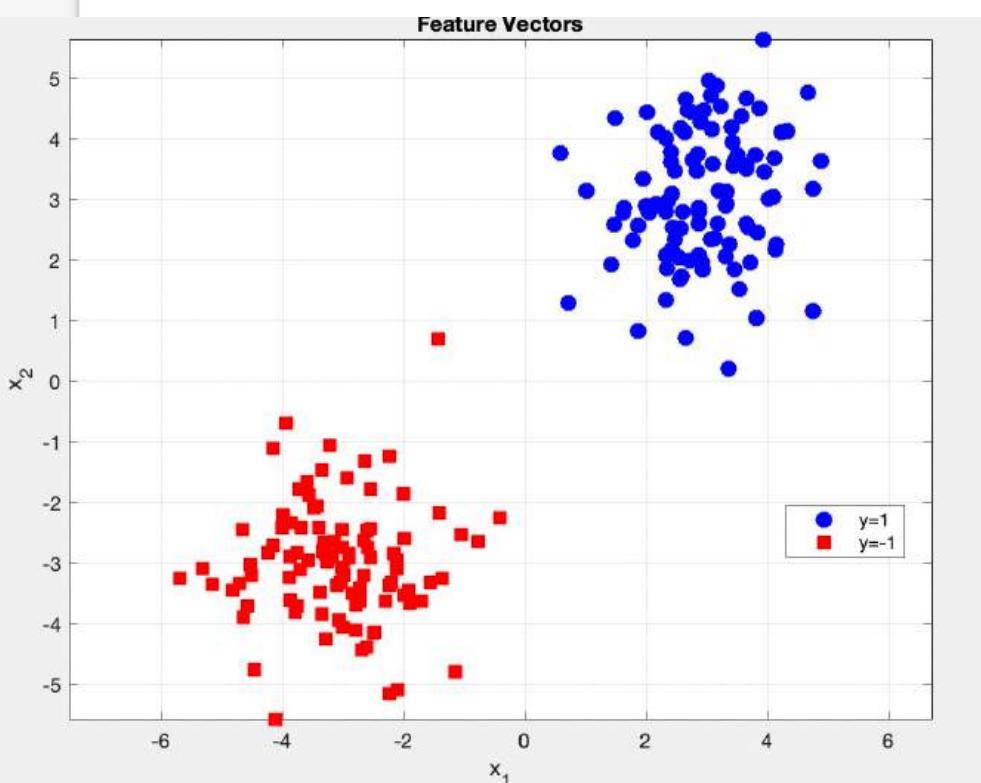
### Command Window

Index	Value
15	2.81583484
17	0.41767449

5) Coding question. We will implement the SVM algorithm with GD to classify two Gaussians in 2D. Dataset is HW7SVM.csv.

(a) In HW7SVM.csv, first 100 rows are the data for cluster 1:  $(x_i, y_i) \in \mathbb{R}^2 \times \mathbb{R}$ ,  $i = 1, \dots, 100$ , with  $y_i = 1$  always. The next 100 rows are the data for cluster 2:  $(x_i, y_i) \in \mathbb{R}^2 \times \mathbb{R}$ ,  $i = 101, \dots, 200$  with  $y_i = -1$  always. Create and turn in a scatter plot of the feature vectors, i.e. the  $x_i$ 's, colored by the label, i.e.,  $y_i$ 's (blue for 1 and red for -1).

```
1 data = readmatrix('HW7SVM.csv');
2 X = data(:, 1:2);
3 y = data(:, 3);
4
5 figure;
6 plot(X(1:100, 1), X(1:100, 2), 'bo', 'MarkerSize', 8, ...
7     'MarkerFaceColor', 'b', 'DisplayName', 'y=1');
8 hold on;
9
10 plot(X(101:200, 1), X(101:200, 2), 'rs', 'MarkerSize', 8, ...
11     'MarkerFaceColor', 'r', 'DisplayName', 'y=-1');
12
13 xlabel('x_1');
14 ylabel('x_2');
15 title('Feature Vectors');
16 legend('Location', 'best');
17 grid on;
18 axis equal;
```



(b) Create a function for the gradient of the loss

$$L(w) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \max(0, 1 - y_i \langle x_i, w \rangle)$$

$$\Delta L(w) = w + \sum_{i=1}^n -y_i x_i \cdot 1_{1-y_i \langle x_i, w \rangle > 0},$$

where  $1_{1-y_i \langle x_i, w \rangle > 0} = \begin{cases} 1 & \text{if } 1 - y_i \langle x_i, w \rangle > 0, \\ 0 & \text{else} \end{cases}$ . Also,  $n=200$ .

To compute the gradient, you'll have to compute an indicator of whether  $1 - y_i \langle x_i, w \rangle > 0$  at every point, and sum up the contribution of this term for all points where its positive.

```
1 function grad = grad(w, X, y)
2     [n, ~] = size(X);
3     w = w(:);
4     grad = w;
5     for i = 1:n
6         dot_product = X(i, :) * w;
7         margin = 1 - y(i) * dot_product;
8         if margin > 0
9             grad = grad - y(i) * X(i, :)';
10        end
11    end
12 end
```



(c) Setting the step size  $\mu = 10^{-4}$  and starting at  $w^{(0)} = (-1, 1)$ , run 1000 iterations of GD. You will create two plots.

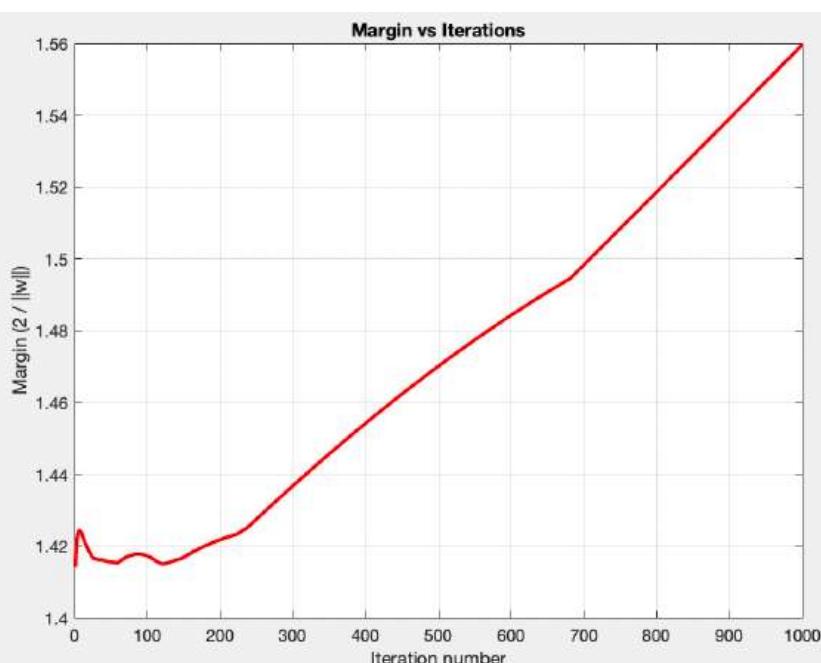
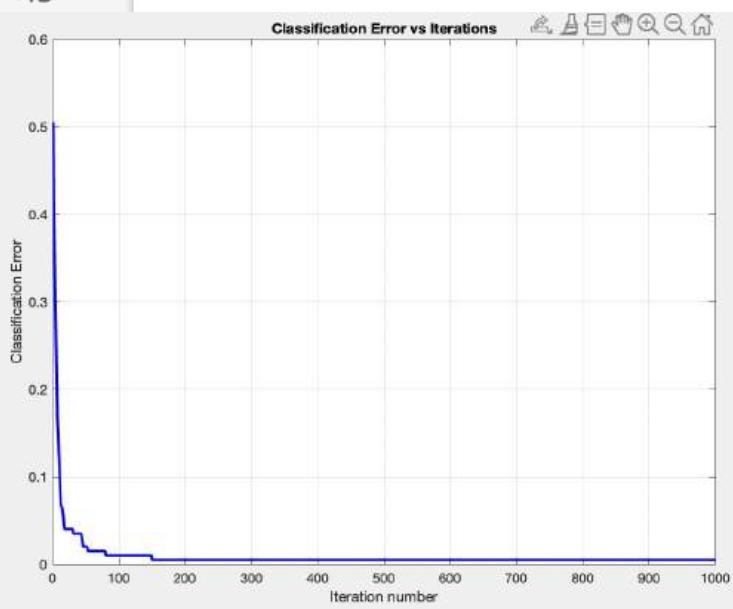
i. Plot the classification error (averaged over all points) as a function of the iterations. The classification of  $x_i$  is determined by  $\text{sign}(\langle x_i, w \rangle)$ .

ii. Plot the margin  $\frac{2}{\|w\|}$  as a function of the iterations. This shows how much of a gap you have between the classes you learned.

```

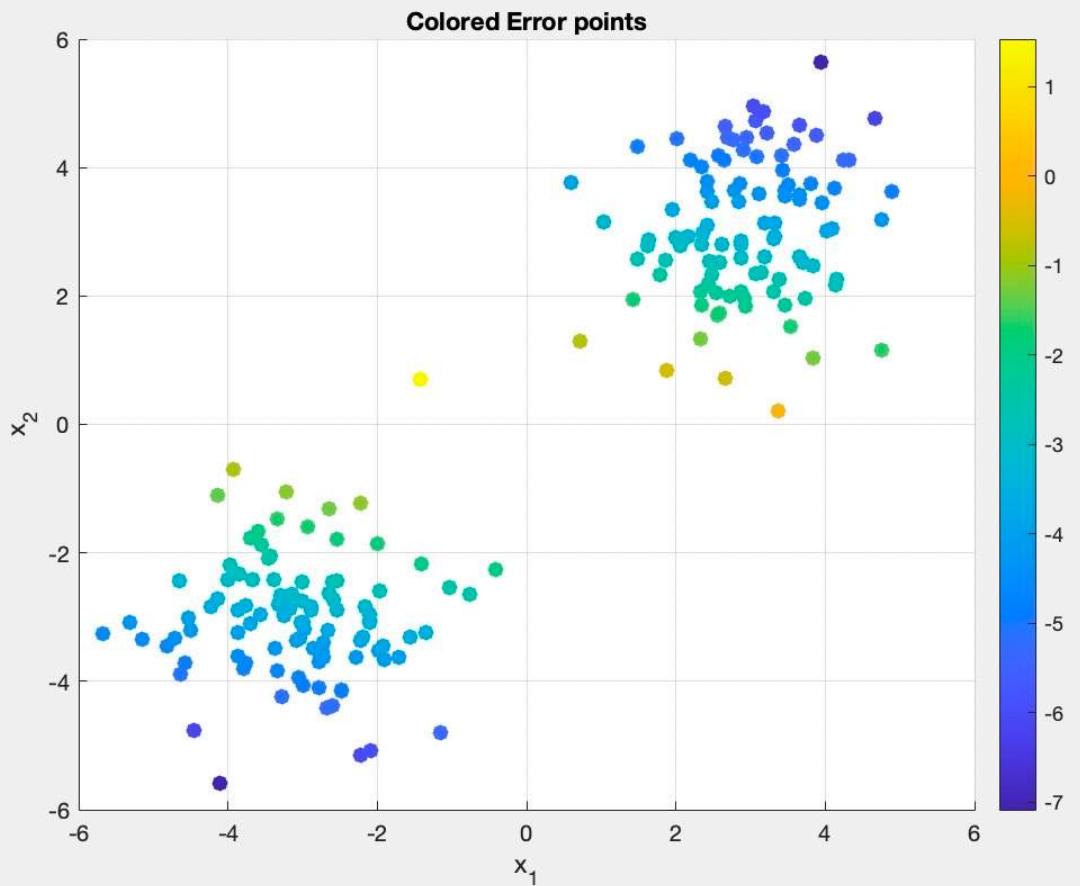
1  data = readmatrix('HW7SVM.csv');
2  X = data(:, 1:2);
3  y = data(:, 3);
4  mu = 1e-4;
5  w0 = [-1; 1];
6  n_iter = 1000;
7  w = w0;
8  error = zeros(n_iter, 1);
9  margin = zeros(n_iter, 1);
10 for iter = 1:n_iter
11     prediction = sign(X * w);
12     misclassified = prediction ~= y;
13     error(iter) = sum(misclassified) / length(y);
14
15     w_norm = norm(w);
16     if w_norm > 0
17         margin(iter) = 2 / w_norm;
18     else
19         margin(iter) = Inf;
20     end
21
22     dot_products = X * w;
23     margin_loss = 1 - y .* dot_products;
24
25     gradient = grad(w, X, y);
26
27     w = w - mu * gradient;
28 end
29
30 figure;
31 plot(1:n_iter, error, 'b-', 'LineWidth', 2);
32 xlabel('Iteration number');
33 ylabel('Classification Error');
34 title('Classification Error vs Iterations');
35 grid on;
36
37 figure;
38 plot(1:n_iter, margin, 'r-', 'LineWidth', 2);
39 xlabel('Iteration number');
40 ylabel('Margin (2 / ||w||)');
41 title('Margin vs Iterations');
42 grid on;
43

```



(d) Create another scatter plot of your data, but this time color the points by the function  $f(x_i) = 1 - y_i \cdot \langle x_i, w \rangle$ . The numbers closest to 0 (positive numbers or largest negative numbers) will show you which points were "most important" in determining the classification.

```
44 f_values = 1 - y .* (X * w);
45 figure;
46 scatter(X(:, 1), X(:, 2), 50, f_values, 'filled');
47 colorbar;
48 xlabel('x_1');
49 ylabel('x_2');
50 title('Colored Error points');
51 grid on;
```



# Homework 6 for Math 173A - Fall 2025

For coding questions, you must submit your code and the requested answers. **Your code must be present to receive points.**

1. Consider the following optimization problem

$$\Phi(x) = \frac{1}{2}x^T \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} x - \sum_{i=1}^2 x_i,$$

where  $x \in \mathbb{R}^2$ . Do the following by hand and clearly write out your steps.

- (a) Find the minimizer of the problem and show that the one you find is indeed the minimizer.  
(b) Is  $\Phi$   $L$ -smooth? Determine the smooth constant  $L$ .  
(c) Perform two steps of the Heavy ball method for  $\Phi$  starting at  $x^{(-1)} = x^{(0)} = (2, 2)$  with stepsize  $\eta = 0.5$  and the parameter  $\beta = 0.1$ . Show your procedures and find what  $x^{(2)}$  is.  
(d) Perform gradient descent with Nesterov's acceleration for  $\Phi$  starting at  $x^{(-1)} = x^{(0)} = (2, 2)$  with stepsize  $\eta = 0.5$  and the parameter  $\beta = 0.1$ . Show your procedures and find what  $x^{(2)}$  is.  
(e) Perform the conjugate gradient descent method (version 0 or 1), where the conjugate directions are initialized and chosen algorithmically, until it finds the minimizer. Show your procedures.
2. **Coding Questions:** Consider Question 4 from Homework 5. We will differentiate 4's and 9's again. This time, use the following two methods to minimize the loss  $F$ :
  - (a) Gradient Descent with momentum
  - (b) Gradient Descent with Nesterov accelerationPlot  $F(w)$  per iteration for each. (you can experiment with the parameters until you find something you like)
3. **Coding Questions:** Let  $A \in \mathbb{R}^{n \times n}$  be a diagonal matrix with diagonal entries  $A_{ii} = i$ , i.e. the entries run from 1 to  $n$ ,

$$A_{ii} = i, \quad \text{i.e. the entries run from 1 to } n,$$

and let  $b \in \mathbb{R}^n$  a vector with all 1 entries. Define the function

$$f(x) = \frac{1}{2}x^T Ax - b^T x.$$

We want to compare the convergence behavior of gradient descent and its accelerated versions . Do the following for  $n = 20$  and  $n = 100$  with initialization  $x^{(0)} = 0$ .

- (a) Find the optimal solution  $x^*$  by solving  $Ax = b$  using a Matlab/Python linear equation solver (or by hand and hard-coding the answer).
- (b) Program and run the gradient descent method for  $f$  with a fixed stepsize. Run the method for 1000 iterations. You may experiment with the stepsize until you see something that works, or use a stepsize dictated by a theorem in the class.
- (c) Program and run the heavy ball method for  $f$ . Run the method for 1000 iterations.
- (d) Program and run gradient descent with Nesterov's acceleration for 1000 iterations.

For items (b) and (c) above, you can experiment with the parameters until you find something you like. Plot the  $f(x^{(t)}) - f(x^*)$  for the three methods in the same figure. In a different figure, plot  $\|x^{(t)} - x^*\|$  for the three methods. If you encounter a number smaller than  $10^{-16}$ , set it to be  $10^{-16}$ . In both plots, make the logarithmic scale for the vertical axis. Comment on the plots.

1) Consider the following optimization problem

$$\Phi(x) = \frac{1}{2} x^T \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} x - \sum_{i=1}^2 x_i ,$$

where  $x \in \mathbb{R}^2$ . Do the following by hand and write out steps.

(a) Find the minimizer of the problem and show that the one you find is indeed the minimizer.

$$\Phi(x) = \frac{1}{2} [x_1 \ x_2] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - (x_1 + x_2)$$

$$= \frac{1}{2} [2x_1 \ x_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - (x_1 + x_2)$$

$$= \frac{1}{2} (2x_1^2 + x_2^2) - (x_1 + x_2)$$

$$= x_1^2 + \frac{1}{2} x_2^2 - x_1 - x_2 \quad 2x_1 = 1 \Rightarrow x_1 = \frac{1}{2}$$

$$\nabla \Phi(x) = \begin{bmatrix} 2x_1 - 1 \\ x_2 - 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \begin{array}{l} x_1 = 1 \\ x_2 = 1 \end{array}$$

$$\Rightarrow x^* = \begin{bmatrix} 1/2 \\ 1 \end{bmatrix}$$

$$\nabla^2 \Phi(x) = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 2-\lambda & 0 \\ 0 & 1-\lambda \end{bmatrix} \Rightarrow \begin{array}{l} \lambda = 2 \\ \lambda = 1 \end{array} \text{ both } \oplus$$

Hessian is PSD  $\Rightarrow \Phi$  is convex, so  $\Phi$  has unique global minimizer  $\Rightarrow x^* = \begin{bmatrix} 1/2 \\ 1 \end{bmatrix}$  is indeed the minimizer

(b) Is  $\Phi$  L-smooth? Determine the smooth constant L.

$$\|\nabla \Phi(x) - \nabla \Phi(y)\| \leq L \|x-y\| \quad \forall x, y \in \mathbb{R}^2$$

$$\begin{aligned}\nabla \Phi(x) &= \begin{bmatrix} 2x_1 - 1 \\ x_2 - 1 \end{bmatrix} & \nabla \Phi(y) &= \begin{bmatrix} 2y_1 - 1 \\ y_2 - 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} & &= \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ A & x & - b & A & y & - b\end{aligned}$$

$$\begin{aligned}\Rightarrow \|\nabla \Phi(x) - \nabla \Phi(y)\| &= \|(Ax-b) - (Ay-b)\| \\ &= \|Ax - Ay\| \\ &= \|A(x-y)\| \\ \Rightarrow \|A(x-y)\| &\leq \lambda_{\max} \|x-y\|\end{aligned}$$

$$\lambda_{\max}(A) = 2 \Rightarrow L = 2$$

Yes,  $\Phi$  is L-smooth with smooth constant  $L=2$ .

(1) Perform 2 steps of the Heavy Ball method for  $\Phi$  starting at  $x^{(1)} = x^{(0)} = (2, 2)$  with stepsize  $\eta = 0.5$  and the parameter  $\beta = 0.1$ . Show your procedures and find what  $x^{(2)}$  is.

$$x^{(t+1)} = x^{(t)} - \mu \nabla f(x^{(t)}) + \beta(x^{(t)} - x^{(t-1)}), \quad t \geq 1$$

$$\begin{aligned}x^{(1)} &= (2, 2) - 0.5 \begin{bmatrix} 2(2) - 1 \\ 2 - 1 \end{bmatrix} + 0.1((2, 2) - (2, 2)) \\ &= (2, 2) - (3/2, 1/2) = (1/2, 3/2)\end{aligned}$$

$$x^{(1)} = (1/2, 3/2)$$

$$\begin{aligned}
 x^{(2)} &= \left(\frac{1}{2}, \frac{3}{2}\right) - 0.5 \begin{bmatrix} 2\left(\frac{1}{2}\right) - 1 \\ \frac{3}{2} - 1 \end{bmatrix} + 0.1 \left(\left(\frac{1}{2}, \frac{3}{2}\right) - (2, 2)\right) \\
 &= \left(\frac{1}{2}, \frac{3}{2}\right) - (0, \frac{1}{4}) + 0.1 \left(-\frac{3}{2}, -\frac{1}{2}\right) \\
 &= (0.5, 1.25) + (-0.15, -0.05) = (0.35, 1.2) \\
 x^{(2)} &= (0.35, 1.2)
 \end{aligned}$$

$$\begin{aligned}
 -\frac{3}{2} &= -1.5 \\
 -\frac{1}{2} &= -0.5
 \end{aligned}$$

(a) Perform gradient descent with Nesterov's acceleration for  $\Phi$  starting at  $x^{(0)} = x^{(1)} = (2, 2)$  with stepsize  $\eta = 0.5$  and the parameter  $\beta = 0.1$ . Show your procedures and find what  $x^{(2)}$  is.

$$y^{(t)} = x^{(t)} + \beta(x^{(t)} - x^{(t-1)})$$

$$x^{(t+1)} = y^{(t)} - \mu \nabla f(y^{(t)})$$

$$y^{(0)} = (2, 2) + 0.1 \cancel{( (2, 2) - (2, 2) ) } = (2, 2)$$

$$x^{(1)} = (2, 2) - (\frac{3}{2}, \frac{1}{2}) = (\frac{1}{2}, \frac{3}{2})$$

$$\begin{aligned}
 y^{(1)} &= \left(\frac{1}{2}, \frac{3}{2}\right) + 0.1 \left(\left(\frac{1}{2}, \frac{3}{2}\right) - (2, 2)\right) \\
 &= (0.5, 1.5) + (-0.15, -0.05) = (0.35, 1.45)
 \end{aligned}$$

$$x^{(2)} = (0.35, 1.45) - 0.5 \begin{bmatrix} 2(0.35) - 1 \\ 1.45 - 1 \end{bmatrix}$$

$$= (0.35, 1.45) - (-0.15, 0.225) = (\frac{1}{2}, 1.225)$$

$$x^{(2)} = (\frac{1}{2}, 1.225)$$

(e) Perform the conjugate gradient descent method (version 0 or 1), where the conjugate directions are initialized and chosen algorithmically, until it finds the minimizer. Show your procedures.

$$x^{(t+1)} = x^{(t)} + \alpha_t p_t, \quad \alpha_t = \arg \min_{\alpha \in \mathbb{R}} \Phi(x^{(t)} + \alpha p_t)$$

$$\text{need } \Phi(x) = \frac{1}{2} x^T A x - b^T x : A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Initialize  $x^{(0)} = (2, 2)$ ,

$$r_0 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix} = \nabla \Phi(x^{(0)})$$

$$p_0 = \begin{bmatrix} -3 \\ -1 \end{bmatrix}, \quad x^{(t+1)} = x^{(t)} + \alpha_t p_t \Rightarrow$$

$$\alpha_0 = \frac{\begin{bmatrix} 3 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \end{bmatrix}}{\begin{bmatrix} -3 & -1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -3 \\ -1 \end{bmatrix}} = \frac{10}{19}$$

$$x^{(1)} = \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \frac{10}{19} \begin{bmatrix} -3 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.421 \\ 1.474 \end{bmatrix}$$

$$r_1 = \begin{bmatrix} 3 \\ 1 \end{bmatrix} + \frac{10}{19} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -3 \\ -1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix} + \begin{bmatrix} -60/19 \\ -10/19 \end{bmatrix} = \begin{bmatrix} -0.158 \\ 0.474 \end{bmatrix}$$

$$\beta_1 = \frac{\begin{bmatrix} -0.158 & 0.474 \end{bmatrix} \begin{bmatrix} -0.158 \\ 0.474 \end{bmatrix}}{\begin{bmatrix} 3 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \end{bmatrix}} = \frac{0.249}{10} = 0.0249$$

$$p_1 = \begin{bmatrix} 0.158 \\ -0.474 \end{bmatrix} + 0.0249 \begin{bmatrix} -3 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.158 \\ -0.474 \end{bmatrix} + \begin{bmatrix} -0.075 \\ -0.0249 \end{bmatrix} = \begin{bmatrix} 0.083 \\ -0.499 \end{bmatrix}$$

$$\alpha_1 = \frac{0.249}{[0.083 \ -0.499] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0.083 \\ -0.499 \end{bmatrix}} = \frac{0.249}{0.263} = 0.95$$

$$x^{(2)} = \begin{bmatrix} 0.421 \\ 1.474 \end{bmatrix} + 0.95 \begin{bmatrix} 0.083 \\ -0.499 \end{bmatrix}$$

$$= \begin{bmatrix} 0.421 \\ 1.474 \end{bmatrix} + \begin{bmatrix} 0.0786 \\ -0.472 \end{bmatrix} = \begin{bmatrix} 0.4996 \\ 1.002 \end{bmatrix} \approx \begin{bmatrix} 1/2 \\ 1 \end{bmatrix}$$

(rounding errors -  $x^{(2)}$  does =  $\begin{bmatrix} 1/2 \\ 1 \end{bmatrix}$  by tm)

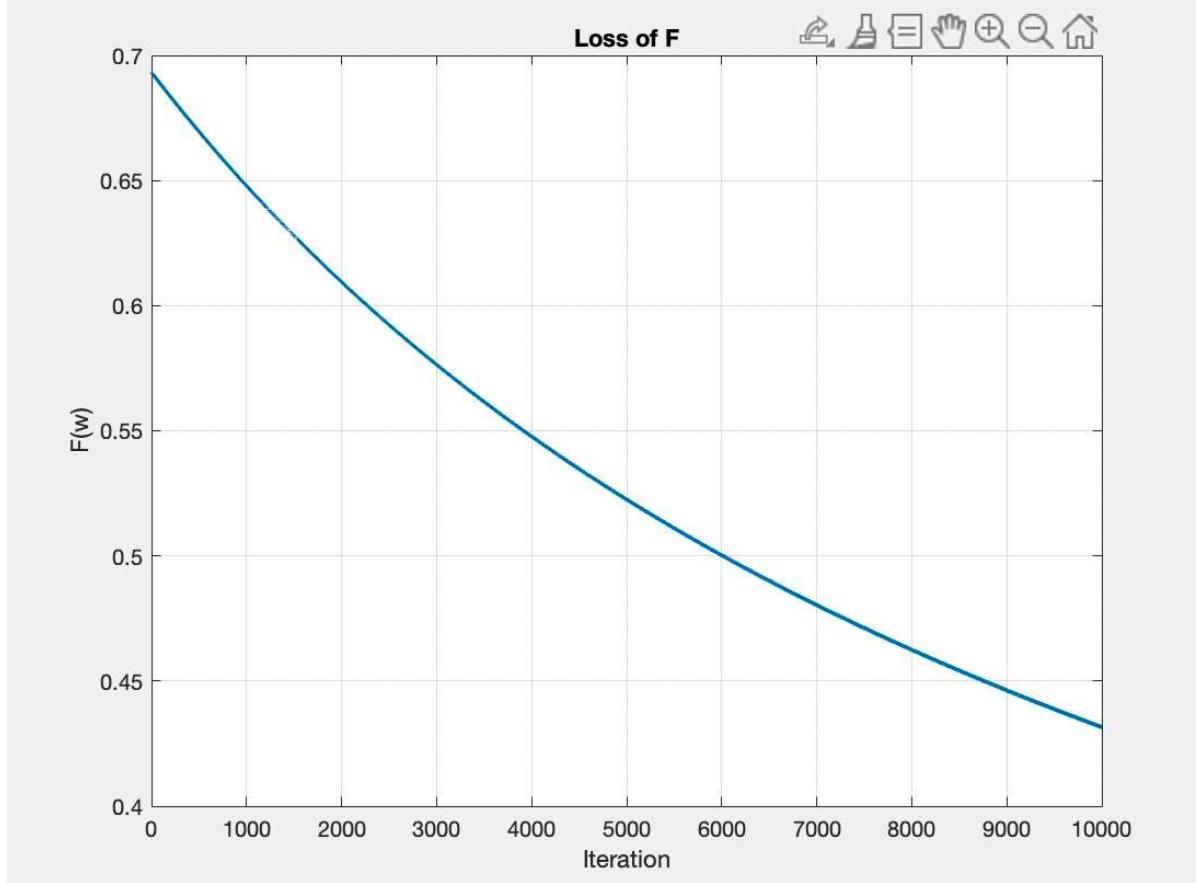
2) Consider Question 4 from Homework 5. We will differentiate  $y$ 's and  $q$ 's again. This time, use the following method to minimize the loss of  $F$ :

### (a) Gradient Descent with momentum (HB)

```

Editor - /Users/sabella/Downloads/gd_momentum.m
form_data.m  mnist_train.csv  mnist_test.csv  gd_momentum.m  gd_NA.m
1 function w = gd_momentum()
2     [x,y] = form_data();
3     [N, d] = size(x);
4     u = 1e-4;
5     max_its = 10000;
6     F_vals = zeros(max_its,1);
7     w = zeros(d,1);
8     beta = 0.1;
9     momentum = zeros(d, 1);
10    for i = 1:max_its
11        inner_prod = x * w;
12        n = length(y);
13        denom = 1 + exp(y .* inner_prod);
14        grad = -(x' * (y ./ denom)) / n;
15        exp_term = exp(-y .* inner_prod);
16        F_vals(i) = (1/N) * sum(log(1 + exp_term));
17
18        momentum = beta * momentum - u * grad;
19        w = w + momentum;
20    end
21
22    figure;
23    plot(1:length(F_vals), F_vals, 'LineWidth', 2);
24    xlabel('Iteration');
25    ylabel('F(w)');
26    title('Loss of F');
27    grid on;
28
29 end

```



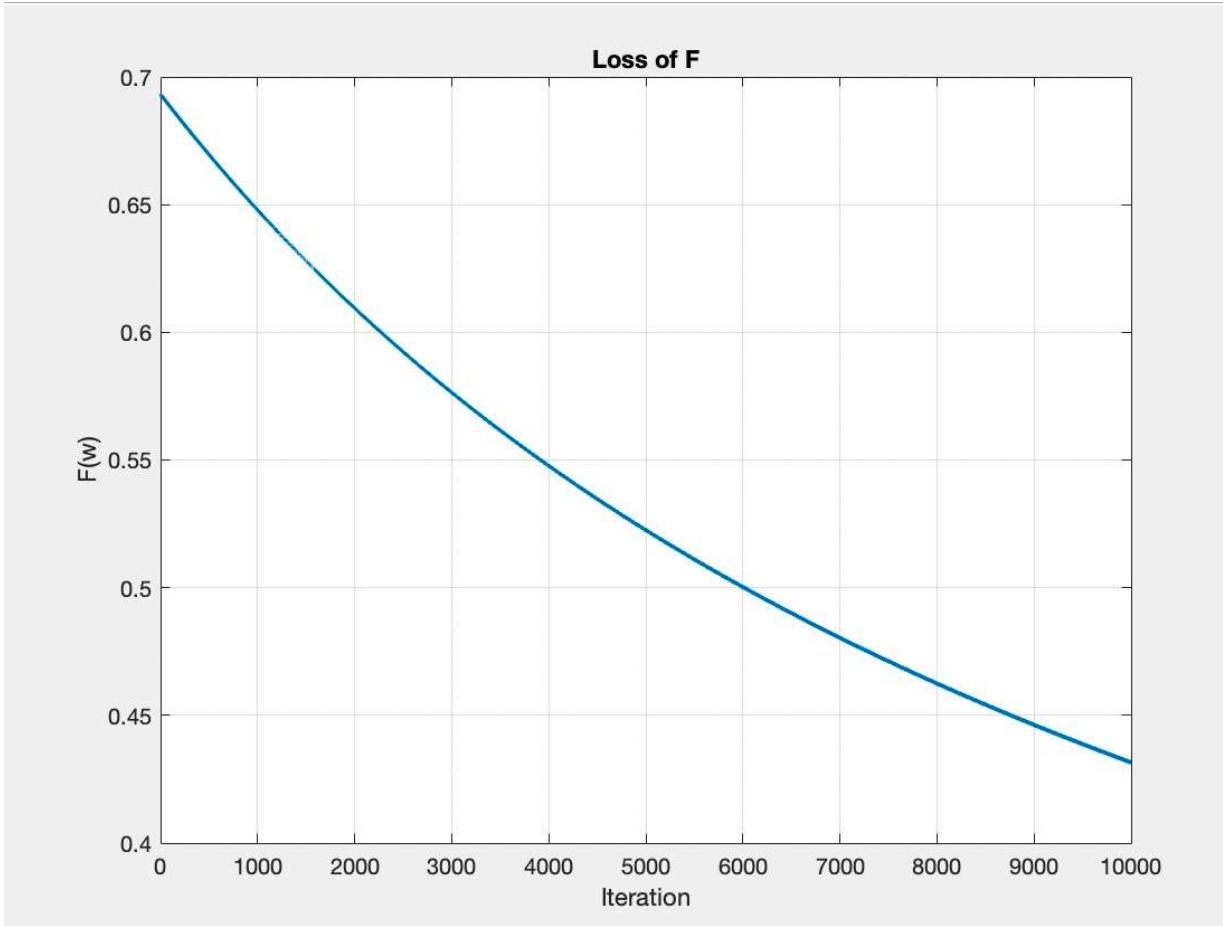
## (b) Gradient Descent with Newton Acceleration

Editor - /Users/sabella/Downloads/gd\_NA.m

```

1 function w = gd_NA()
2     [x,y] = form_data();
3     [N, d] = size(x);
4     u = 1e-4;
5     max_its = 10000;
6     F_vals = zeros(max_its,1);
7     w = zeros(d,1);
8     beta = 0.1;
9     w_prev = w;
10    for i = 1:max_its
11        w_y = w + beta * (w - w_prev);
12        inner_prod = x * w_y;
13        n = length(y);
14        denom = 1 + exp(y .* inner_prod);
15        grad_lookahead = -(x' * (y ./ denom)) / n;
16        inner_prod_current = x * w;
17        exp_term_current = exp(-y .* inner_prod_current);
18        F_vals(i) = (1/N) * sum(log(1 + exp_term_current));
19
20        w_prev = w;
21        w = w_y - u * grad_lookahead;
22    end
23
24    figure;
25    plot(1:length(F_vals), F_vals, 'LineWidth', 2);
26    xlabel('Iteration');
27    ylabel('F(w)');
28    title('Loss of F');
29    grid on;
30
end

```



3) Let  $A \in \mathbb{R}^{n \times n}$  be a diagonal matrix with diagonal entries  $A_{ii} = i$ , i.e. the entries run from 1 to  $n$ , and let  $b \in \mathbb{R}^n$  be a vector with all 1 entries. Define

$$f(x) = \frac{1}{2} x^T A x - b^T x.$$

We want to compare the convergence of gradient descent and its accelerated versions. Do the following for  $n=20$  and  $n=100$  with initialization  $x^{(0)} = 0$ .

(a) Find the optimal solution  $x^*$  by solving  $Ax=b$  using a Matlab/Python linear equation solver (or by hand and hard-coding the answer).

```
Editor - /Users/sabella/Downloads/untitled4.m
form_data.m mnist_train.csv mnist_test.csv gd_momentum.m gd_NA.m untitled4.m +
1 n1 = 20;
2 A1 = diag(1:n1);
3 b1 = ones(n1,1);
4 x_star_1 = A1 \ b1;
5
6 n2 = 100;
7 A2 = diag(1:n2);
8 b2 = ones(n2,1);
9 x_star_2 = A2 \ b2;
10
11 fprintf('x* for n = 20: \n')
12 disp(x_star_1)
13 fprintf('x* for n = 100: \n')
14 disp(x_star_2)|
```

## Command Window

```
>> untitled4
*x for n = 20:
 1.0000
 0.5000
 0.3333
 0.2500
 0.2000
 0.1667
 0.1429
 0.1250
 0.1111
 0.1000
 0.0909
 0.0833
 0.0769
 0.0714
 0.0667
 0.0625
 0.0588
 0.0556
 0.0526
 0.0500
```

x\* for n = 100:

```
 1.0000
 0.5000
 0.3333
 0.2500
 0.2000
 0.1667
 0.1429
 0.1250
 0.1111
 0.1000
 0.0909
 0.0833
 0.0769
 0.0714
 0.0667
 0.0625
 0.0588
 0.0556
 0.0526
 0.0500
 0.0476
 0.0455
 0.0435
 0.0417
 0.0400
 0.0385
 0.0370
 0.0357
 0.0345
 0.0333
 0.0323
 0.0312
```

(please forgive the  
strange layout -  
trying to save space)

(b) Program and run the GD method for f with a fixed stepsize. Run the method for 1000 iterations. You may experiment with the stepsize until you see something that works, or use a stepsize dictated by a theorem in class.

Editor - /Users/sabella/Downloads/gd\_hw6.m

```
+4 gd_NA.m * untitled4.m * gd_hw6.m * untitled6 * * untitled7 * * untitled8 * +
1 function x = gd_hw6()
2     n = 20;
3     max_its = 1000;
4     A = diag(1:n);
5     b = ones(n, 1);
6     x = zeros(n, 1);
7     f_vals = zeros(max_its, 1);
8     L = n;
9     u = 1 / L;
10    for i = 1:max_its
11        grad = A * x - b;
12        f_vals(i) = 0.5 * x' * A * x - b' * x;
13        x = x - u * grad;
14    end
15    disp(x);
16 end
```

&gt;&gt; gd\_hw6 (n=20)

```
1.0000
0.5000
0.3333
0.2500
0.2000
0.1667
0.1429
0.1250
0.1111
0.1000
0.0909
0.0833
0.0769
0.0714
0.0667
0.0625
0.0588
0.0556
0.0526
0.0500
```

(n=100)

```
>> gd_hw6 (n=100)
 1.0000
 0.5000
 0.3333
 0.2500
 0.2000
 0.1667
 0.1429
 0.1250
 0.1111
 0.1000
 0.0909
 0.0833
 0.0769
 0.0714
 0.0667
 0.0625
 0.0588
 0.0556
 0.0526
 0.0500
 0.0476
 0.0455
 0.0435
 0.0417
 0.0400
 0.0385
 0.0370
 0.0357
 0.0345
 0.0333
 0.0323
 0.0312
```

(c) Program and run the HB method for f. Run the method for 1000 iterations.

```
Editor - /Users/sabella/Downloads/gd_HB.m *
+2 mnist_test.csv x untitled4.m x gd_hw6.m x gd_HB.m * x untitled7 * x untitled8 * x +
1 function x = gd_HB()
2     n = 20;
3     max_its = 1000;
4     A = diag(1:n);
5     b = ones(n, 1);
6     x = zeros(n, 1);
7     f_vals = zeros(max_its, 1);
8     L = n;
9     u = 1 / L;
10    beta = 0.1;
11    momentum = zeros(n, 1);
12    for i = 1:max_its
13        grad = A * x - b;
14        f_vals(i) = 0.5 * x' * A * x - b' * x;
15        momentum = beta * momentum - u * grad;
16        x = x + momentum;
17    end
18    disp(x);
19 end
```

>> gd\_HB (n = 20)

```
1.0000
0.5000
0.3333
0.2500
0.2000
0.1667
0.1429
0.1250
0.1111
0.1000
0.0909
0.0833
0.1000
0.0909
0.0833
0.0769
0.0714
0.0667
0.0625
0.0588
0.0556
0.0667
0.0625
0.0588
0.0556
0.0526
0.0500
0.0476
0.0455
0.0435
0.0417
0.0400
0.0385
0.0370
0.0357
0.0345
0.0333
0.0323
0.0312
0.0303
```

>> gd\_HB (n = 100)

1.0000	0.0294
0.5000	0.0286
0.3333	0.0278
0.2500	0.0270
0.2000	0.0263
0.1667	0.0256
0.1429	0.0250
0.1250	0.0244
0.1111	0.0238
0.1000	0.0233
0.0909	0.0227
0.0833	0.0222
0.1000	0.0217
0.0769	0.0213
0.0714	0.0208
0.0667	0.0204
0.0625	0.0200
0.0588	0.0196
0.0556	0.0192
0.0667	0.0189
0.0625	0.0185
0.0588	0.0182
0.0556	0.0179
0.0526	0.0175
0.0500	0.0172
0.0476	0.0169
0.0455	0.0167
0.0435	0.0164
0.0417	0.0161
0.0400	0.0159
0.0385	0.0156
0.0370	0.0154
0.0357	0.0152
0.0345	0.0149

(d) Program and run GD with NA for 1000 iterations.

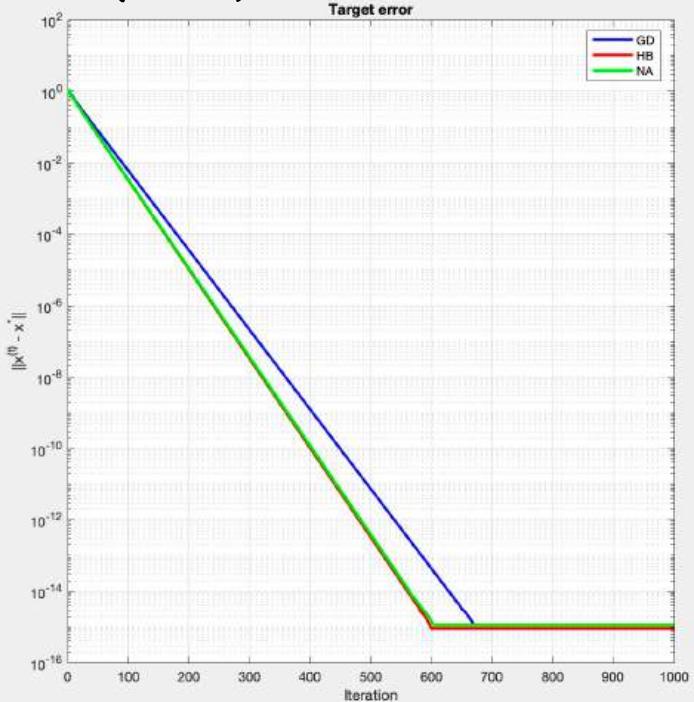
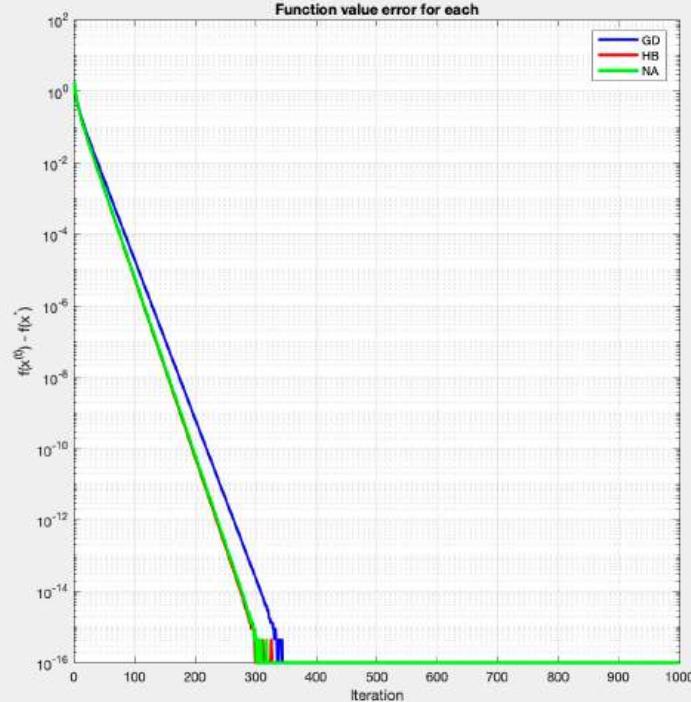
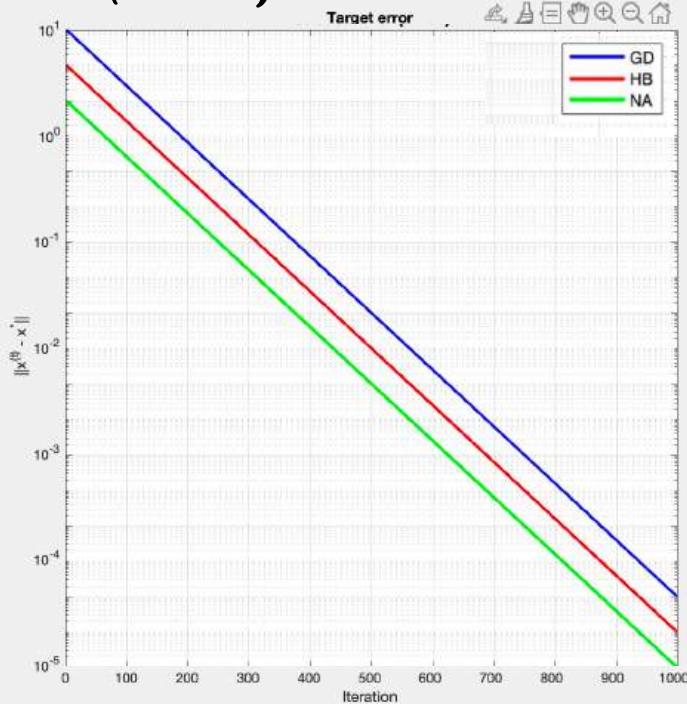
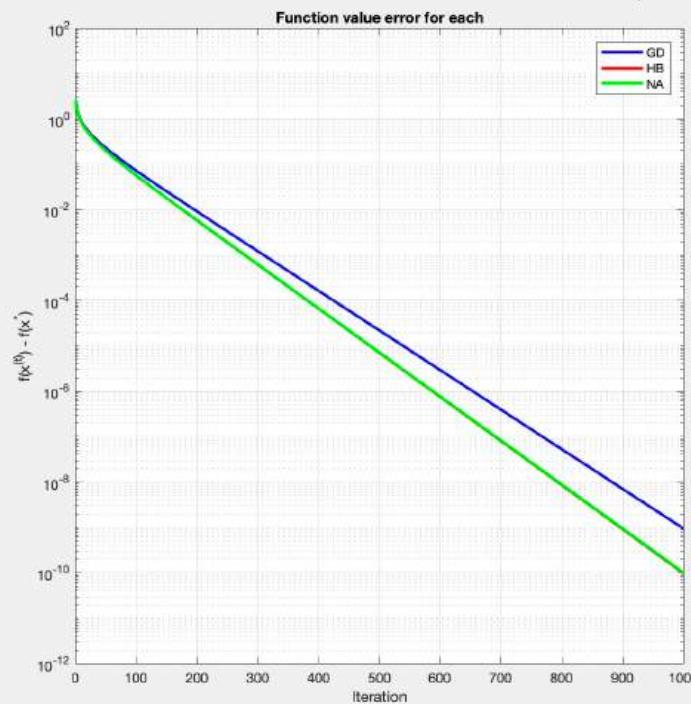
Editor - /Users/sabella/Downloads/gd\_NA\_2.m

```
1 function x = gd_NA_2()
2     n = 20;
3     max_its = 1000;
4     A = diag(1:n);
5     b = ones(n, 1);
6     x = zeros(n, 1);
7     f_vals = zeros(max_its, 1);
8     L = n;
9     u = 1 / L;
10    beta = 0.1;
11    x_prev = x;
12    for i = 1:max_its
13        y = x + beta*(x-x_prev);
14        grad = A * y - b;
15        f_vals(i) = 0.5 * x' * A * x - b' * x;
16
17        x_prev = x;
18        x = y - u*grad;
19    end
20    disp(x);|
21 end
```

>> gd\_NA\_2 (n = 20)

x <sup>(t)</sup>	f(x <sup>(t)</sup> )
1.0000	0.0294
0.5000	0.0286
0.3333	0.0278
0.2500	0.0270
0.2000	0.0263
0.1667	0.0256
0.1429	0.0250
0.1250	0.0244
0.1111	0.0238
0.1000	0.0233
0.0909	0.0227
0.0833	0.0222
0.0769	0.0217
0.0714	0.0213
0.0667	0.0208
0.0625	0.0204
0.0588	0.0200
0.0556	0.0196
0.0526	0.0192
0.0500	0.0189
0.0476	0.0185
0.0455	0.0182
0.0435	0.0179
0.0417	0.0175
0.0400	0.0172
0.0385	0.0169
0.0370	0.0167
0.0357	0.0164
0.0345	0.0161
0.0333	0.0159
0.0323	0.0156
0.0312	0.0154
0.0303	0.0152
	0.0149

For items (b) and (c) above, you can experiment with the parameters until you find something you like. Plot the  $f(x^{(t)}) - f(x^*)$  for the three methods in the same figure. In a different figure, plot  $\|x^{(t)} - x^*\|$  for the three methods. If you encounter a number smaller than  $10^{-16}$ , set it to be  $10^{-16}$ . In both plots, make the logarithmic scale for the vertical axis. Comment on the plots.

Comparison of the three methods ( $n=20$ )Comparison of the three methods ( $n=100$ )

for  $n=20$ , we see that the regular GD has the highest error during the iterate, for both  $f(x^{(t)})$  and  $x^{(t)}$ , and thus rounds out to about the same as HB and NA.

For  $n=100$ , the actual value of the target error  $\|x^{(t)} - x^*\|$  has a clear spread with NA producing the least error, regular GD producing the most, and HB exactly in the middle. Again, for the function value, NA and HB are almost exactly the same,

while regular GD produces slightly more error.

```
Editor - /Users/sabella/Downloads/plot_three_methods.m*
form_data.m  x  mnist_train.csv  x  mnist_test.csv  x  untitled4.m  x  gd_hw6.m  x  gd_HB.m  :
1 function plot_three_methods()
2     n = 100;
3     max_its = 1000;
4
5     A = diag(1:n);
6     b = ones(n, 1);
7     x_star = A \ b;
8     f_star = 0.5 * x_star' * A * x_star - b' * x_star;
9
10    gd_hw6();
11    gd_HB();
12    gd_NA_2();
13
14    gd_data = load('gd_hw6.mat');
15    hb_data = load('gd_HB.mat');
16    na_data = load('gd_NA_2.mat');
17
18    f_gd = gd_data.f_history;
19    f_hb = hb_data.f_history;
20    f_na = na_data.f_history;
21
22    f_gap_gd = max(f_gd - f_star, 1e-16);
23    f_gap_hb = max(f_hb - f_star, 1e-16);
24    f_gap_na = max(f_na - f_star, 1e-16);
25
26    error_gd = zeros(max_its, 1);
27    error_hb = zeros(max_its, 1);
28    error_na = zeros(max_its, 1);
29
30    for i = 1:max_its
31        error_gd(i) = max(norm(gd_data.x_history(i, :) - x_star), 1e-16);
32        error_hb(i) = max(norm(hb_data.x_history(i, :) - x_star), 1e-16);
33        error_na(i) = max(norm(na_data.x_history(i, :) - x_star), 1e-16);
34    end
35
36    figure('Position', [100, 100, 1200, 500]);
37
38    subplot(1, 2, 1);
39    semilogy(1:length(f_gap_gd), f_gap_gd, 'b-', 'LineWidth', 2); hold on;
40    semilogy(1:length(f_gap_hb), f_gap_hb, 'r-', 'LineWidth', 2);
41    semilogy(1:length(f_gap_na), f_gap_na, 'g-', 'LineWidth', 2);
42    xlabel('Iteration');
43    ylabel('f(x^{(t)}) - f(x^*)');
44    title(sprintf('Function value error for each'));
45    legend('GD', 'HB', 'NA', 'Location', 'best');
46    grid on;
47
48    subplot(1, 2, 2);
49    semilogy(1:length(error_gd), error_gd, 'b-', 'LineWidth', 2); hold on;
50    semilogy(1:length(error_hb), error_hb, 'r-', 'LineWidth', 2);
51    semilogy(1:length(error_na), error_na, 'g-', 'LineWidth', 2);
52    xlabel('Iteration');
53    ylabel('||x^{(t)} - x^*||');
54    title(sprintf('Target error'));
55    legend('GD', 'HB', 'NA', 'Location', 'best');
56    grid on;
57
58    sgttitle(sprintf('Comparison of the three methods'));
59
60 end
```

# Homework 5 for Math 173A - Fall 2025

For coding questions, you must submit your code and the requested answers. **Your code must be present to receive points.**

1. Consider the function  $f(x_1, x_2) = (2x_1 - 1)^4 + (x_1 + x_2 - 1)^2$ . Recall the minimizer is  $(\frac{1}{2}, \frac{1}{2})$  for this problem.
  - (a) Starting at  $x^{(0)} = (0, 0)$ , perform two steps of Newton's method to find  $x^{(2)}$ . Show your work.
  - (b) Can you apply the theorem of convergence for Newton's method learned in the lecture to this problem, even if the initial  $x^{(0)}$  is very close to  $(\frac{1}{2}, \frac{1}{2})$ ? Hint: check whether the Hessian at  $x^*$  is invertible.
  - (c) Consider the function  $g(x_1, x_2) = (2x_1 - 1)^2 + (x_1 + x_2 - 1)^2$ . Redo the subquestion (a) for this function  $g$ .

**Bonus question:** If you start Newton's method at  $x^{(0)} = (0, 0)$  for the function  $f$ , does Newton's method converge to  $x^*$ ? Justify your answer.

2. **Coding Question:** We'll consider the function  $f(x_1, x_2) = 200(x_2 - x_1^2)^2 + (1 - x_1)^2$ .
  - (a) Write a computer program to run Newton's method on this problem.
  - (b) Write a computer program to run Gradient Descent with a fixed step size  $\mu = 10^{-3}$  on this problem.
  - (c) Write a computer program to run Gradient Descent with backtracking line-search for this problem (you may set  $\beta$  and  $\gamma$  as you wish).
  - (d) Starting at the same  $x^{(0)}$ , run each of the four algorithms and plot  $\|x^{(t)} - x^*\|$  for each, in the same figure. In a separate figure, plot  $f(x^{(t)}) - f(x^*)$  for each of them. Comment on the performance. Note that the minimizer  $x^* = (1, 1)$  for this problem.
3. **Coding Question:** We will redo the MNIST coding question but **for differentiating 4's and 9's**. You can reuse the template for this problem.
  - (a) Repeat parts (b), (c), and (e) from Q4 of HW4 for digits of 4s and 9s. Comment on the difference between the results and propose a reason as to why the performance did or did not change.

- (b) Implement and run gradient descent with **backtracking line search**. At each iteration  $t$ , initialize the step size  $\mu = 10^{-1}$ , and use  $\gamma = 0.5$  and  $\beta = 0.8$  to determine the correct  $\mu^{(t)}$ . Run your algorithm for at least 10,000 iterations and plot the loss curve  $F(w^{(t)})$  as a function of  $t$ .
4. **Coding Question:** Again, we will redo the MNIST coding question from HW4 but using different versions of gradient descent

$$w^{(t+1)} = w^{(t)} - \mu p^{(t)}.$$

This time, we will **differentiate 4's and 9's**. That is, instead of classifying images of 0 and 1, you classify the images of 4 and 9. You can reuse the template from the previous homework for loading / formatting MNIST.

- (a) Implement and run  $L^\infty$  gradient descent with step size  $\mu = 10^{-8}$ . Run your algorithm for at least  $10^4$  iterations and initialize with  $w^{(0)} = 0$  (i.e. the zero vector). **Recall:**  $L^\infty$  gradient descent uses steps

$$p^{(t)} = \text{sign}(\nabla F(w^{(t)})) \|\nabla F(w^{(t)})\|_1.$$

- (b) Implement and run  $L^1$  gradient descent (aka. coordinate descent) with step size  $\mu = 10^{-4}$ . Run your algorithm for at least  $10^4$  iterations and initialize with  $w^{(0)} = 0$  (i.e. the zero vector). **Recall:**  $L^1$  gradient descent uses steps

$$p^{(t)} = \text{sign}(\nabla_{j^*} F(w^{(t)})) \|\nabla F(w^{(t)})\|_\infty e_{j^*},$$

where  $j^*$  is the location of the largest entry of the gradient, and  $e_j$  is the zero vector with a 1 in the  $j^{th}$  entry.

- (c) Compare these two descent plots of  $F(w)$ , along with the analogous plot for gradient descent from HW4. Which performs best, and do you have an argument for why? Do you think the performance would change with different step sizes?
- (d) For the coordinate descent problem, rerun gradient descent but store a running sum of which entry of  $p^{(t)}$  is nonzero at each iteration (not the actual value of the direction vector, just  $e_{j^*}$ ). This will result in a size 784 vector of mostly zeros, and should have integers at various entries whose sum equals the number of iterations. Reshape this vector to be a 28x28 image and display the result. Why do you think these are the pixels that were chosen in the gradient? How can you use this to interpret the algorithm and its results?

1) Consider the function  $f(x_1, x_2) = (2x_1 - 1)^4 + (x_1 + x_2 - 1)^2$ . Recall the minimizer is  $(\frac{1}{2}, \frac{1}{2})$  for this problem.

(a) Starting at  $x^{(0)} = (0, 0)$ , perform two steps of Newton's method to find  $x^{(2)}$ .

$$x^{(t+1)} = x^{(t)} - [\nabla^2 f(x^{(t)})]^{-1} \nabla f(x^{(t)})$$

$$\nabla f(x_1, x_2) = (8(2x_1 - 1)^3 + 2(x_1 + x_2 - 1), 2(x_1 + x_2 - 1))^T$$

$$\nabla^2 f(x_1, x_2) = \begin{bmatrix} 48(2x_1 - 1)^2 + 2 & 2 \\ 2 & 2 \end{bmatrix}$$

$$\nabla f(0, 0) = (8(0 - 1)^3 + 2(0 + 0 - 1), 2(0 + 0 - 1)) = (-10, -2)^T$$

$$\nabla^2 f(0, 0) = \begin{bmatrix} 48(0 - 1)^2 + 2 & 2 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 50 & 2 \\ 2 & 2 \end{bmatrix}$$

$$(\nabla^2 f(0, 0))^{-1} = \frac{1}{100 - 4} \begin{bmatrix} 2 & -2 \\ -2 & 50 \end{bmatrix} = \frac{1}{96} \begin{bmatrix} 2 & -2 \\ -2 & 50 \end{bmatrix}$$

$$x^{(1)} = (0, 0)^T - \frac{1}{96} \begin{bmatrix} 2 & -2 \\ -2 & 50 \end{bmatrix} \begin{bmatrix} -10 \\ -2 \end{bmatrix}$$

$$= (0, 0)^T - \frac{1}{96} \begin{bmatrix} -16 \\ -80 \end{bmatrix} = \frac{1}{96} \begin{bmatrix} 16 \\ 80 \end{bmatrix} = \begin{bmatrix} 1/6 \\ 5/6 \end{bmatrix}$$

$$\begin{aligned} \nabla f(\frac{1}{6}, \frac{5}{6}) &= (8(2(\frac{1}{6}) - 1)^3 + 2(\frac{1}{6} + \frac{5}{6} - 1), 2(\frac{1}{6} + \frac{5}{6} - 1))^T \\ &= (8(-\frac{2}{3})^3, 0) = \begin{bmatrix} -64/27 \\ 0 \end{bmatrix} \end{aligned}$$

$$\nabla^2 f(1/6, 5/6) = \begin{bmatrix} 48(2(1/6)-1)^2+2 & 2 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 48(4/9)+2 & 2 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 70/3 & 2 \\ 2 & 2 \end{bmatrix}$$

$$(\nabla^2 f(1/6, 5/6))^{-1} = \frac{1}{140/3 - 4} \begin{bmatrix} 2 & -2 \\ -2 & 70/3 \end{bmatrix} = \frac{3}{128} \begin{bmatrix} 2 & -2 \\ -2 & 70/3 \end{bmatrix}$$

$$x^{(2)} = \begin{bmatrix} 1/6 \\ 5/6 \end{bmatrix} - \frac{3}{128} \begin{bmatrix} 2 & -2 \\ -2 & 70/3 \end{bmatrix} \begin{bmatrix} -64/27 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1/6 \\ 5/6 \end{bmatrix} - \frac{3}{128} \begin{bmatrix} -128/27 \\ 128/27 \end{bmatrix} = \begin{bmatrix} 1/6 \\ 5/6 \end{bmatrix} - \begin{bmatrix} -1/9 \\ 1/9 \end{bmatrix} = \begin{bmatrix} 5/18 \\ 13/18 \end{bmatrix}$$

(b) Can you apply the theorem of convergence for Newton's method learned in lecture to this problem, even if the initial  $x^{(0)}$  is very close to  $(1/2, 1/2)$ ? Hint: check whether Hessian at  $x^*$  is invertible.

$$\nabla^2 f(x^*) = \nabla^2 f(1/2, 1/2) = \begin{bmatrix} 48(2(1/2)-1)^2+2 & 2 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

$$(\nabla^2 f(x^*))^{-1} = \frac{1}{4-4} \begin{bmatrix} 2 & -2 \\ -2 & 2 \end{bmatrix} = \text{not invertible},$$

So no, we cannot apply the theorem of convergence for Newton's method, even if initial  $x^{(0)}$  is very close to  $(1/2, 1/2)$ .

(c) Consider the function  $g(x_1, x_2) = (2x_1 - 1)^2 + (x_1 + x_2 - 1)^2$ . Redo (a) for this function  $g$ .

$$x^{(t+1)} = x^{(t)} - [\nabla^2 g(x^{(t)})]^{-1} \nabla g(x^{(t)})$$

$$\nabla g(x_1, x_2) = (10x_1 + 2x_2 - 6, 2x_1 + 2x_2 - 2)^T$$

$$\nabla^2 g(x_1, x_2) = \begin{bmatrix} 10 & 2 \\ 2 & 2 \end{bmatrix}$$

$$\nabla g(0,0) = (-6, -2)^T \quad \nabla^2 g(0,0) = \begin{bmatrix} 10 & 2 \\ 2 & 2 \end{bmatrix}$$

$$(\nabla^2 g(0,0))^{-1} = \frac{1}{20-4} \begin{bmatrix} 2 & -2 \\ -2 & 10 \end{bmatrix} = \frac{1}{16} \begin{bmatrix} 2 & -2 \\ -2 & 10 \end{bmatrix}$$

$$x^{(1)} = (0,0) - \frac{1}{16} \begin{bmatrix} 2 & -2 \\ -2 & 10 \end{bmatrix} \begin{bmatrix} -6 \\ -2 \end{bmatrix} = -\frac{1}{16} \begin{bmatrix} -8 \\ -8 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix} = x^*$$

2) Consider the function  $f(x_1, x_2) = 200(x_2 - x_1^2)^2 + (1 - x_1)^2$

(a)

The screenshot shows a MATLAB code editor window titled "Editor - /Users/sabella/Downloads/Newton.m \*". The code is a function named "Newton" that implements the Newton-Raphson method to find the minimum of the given function. It uses a while loop to iteratively update the solution vector  $x$  until the norm of the gradient is less than a tolerance  $tol$ . The Hessian matrix  $H$  is calculated at each iteration, and the search direction  $d$  is determined by solving the linear system  $-H \backslash g$ .

```
function x_min = Newton(x0, tol, max_its)
    x = x0;
    i = 0;

    while i < max_its
        x1 = x(1);
        x2 = x(2);
        g1 = -400 * x1 * (x2 - x1^2) - 2 * (1 - x1);
        g2 = 400 * (x2 - x1^2);
        g = [g1; g2];
        if norm(g) < tol
            break;
        end
        H11 = 1200 * x1^2 - 400 * x2 + 2;
        H12 = -400 * x1;
        H21 = H12;
        H22 = 400;
        H = [H11, H12; H21, H22];

        d = -H \ g;
        x = x + d;
        i = i + 1;
    end
    x_min = x;
end
```

(b)

The screenshot shows a MATLAB code editor window titled "Editor - /Users/sabella/Downloads/GradDesc2.m". This script, named "GradDesc2", implements gradient descent to find the minimum of the same function. It uses a while loop to iteratively update the solution vector  $x$  by subtracting a step size  $u$  times the negative gradient. The loop continues until the norm of the gradient is less than a tolerance  $tol$ .

```
function x_min = GradDesc2(x0, tol, max_its)
    x = x0;
    i = 0;
    u = 0.001;

    while i < max_its
        x1 = x(1);
        x2 = x(2);
        g1 = -400 * x1 * (x2 - x1^2) - 2 * (1 - x1);
        g2 = 400 * (x2 - x1^2);
        g = [g1; g2];
        if norm(g) < tol
            break;
        end
        x = x - u*g;
        i = i + 1;
    end
    x_min = x;
end
```

(c)

```

Editor - /Users/sabella/Downloads/gd_backtrack.m
Newton.m  GradDesc2.m  gd_backtrack.m  +
1 function x_min = gd_backtrack(x0, max_its)
2     beta = 0.5;
3     gamma = 0.1;
4     x = x0;
5     for i = 1:max_its
6         x1 = x(1);
7         x2 = x(2);
8         g1 = -400 * x1 * (x2 - x1^2) - 2 * (1 - x1);
9         g2 = 400 * (x2 - x1^2);
10        g = [g1; g2];
11        alpha = 1.0;
12        f_curr = 200*(x2-x1^2)^2 + (1-x1)^2;
13
14        while true
15            x_new = x - alpha * g;
16            x1_new = x_new(1);
17            x2_new = x_new(2);
18            f_new = 200*(x2_new-x1_new^2)^2 + (1-x1_new)^2;
19            if f_new <= f_curr - gamma * alpha * (g' * g)
20                break;
21            end
22            alpha = beta * alpha;
23        end
24        x = x - alpha * g;
25    end
26    x_min = x;
27 end

```

(d) Starting at the same  $x^{(0)}$ , run each of the four algos and plot  $\|x^{(t)} - x^*\|$  for each, in the same figure. In a separate figure, plot  $f(x^{(t)}) - f(x^*)$  for each of them. Comment on performance. Note that minimizer  $x^* = (1, 1)$  for this problem.

```

clear; clc;
x_star = [1; 1];
x0 = [0; 0];
max_its = 100;

d_newton = zeros(1, max_its);
d_GradDesc2 = zeros(1, max_its);
d_gd_backtrack = zeros(1, max_its);

newton_vals = zeros(1, max_its);
GradDesc2_vals = zeros(1, max_its);
gd_backtrack_vals = zeros(1, max_its);

x = x0;
for i = 1:max_its
    x = Newton(x, 1e-6, 1);
    d_newton(i) = norm(x - x_star);
    x1 = x(1);
    x2 = x(2);
    newton_vals(i) = 200*(x2-x1^2)^2 + (1-x1)^2;
end

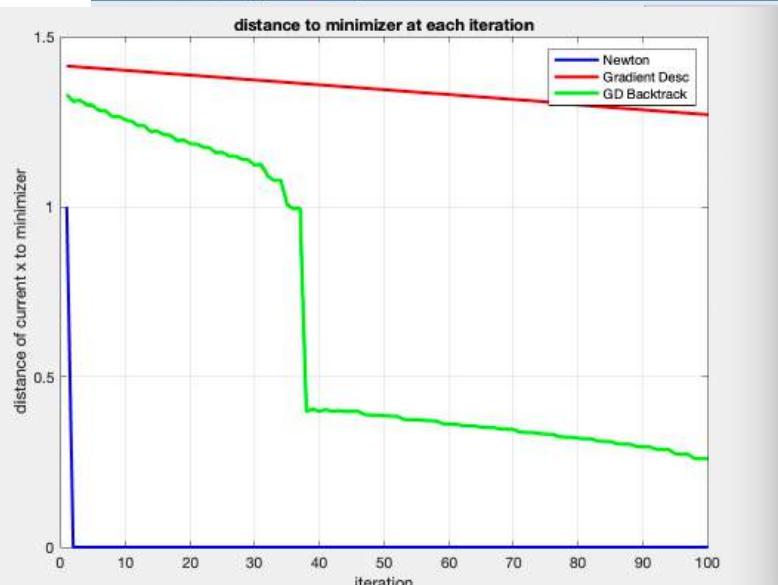
x = x0;
for i = 1:max_its
    x = GradDesc2(x, 1e-6, 1);
    d_GradDesc2(i) = norm(x - x_star);
    x1 = x(1);
    x2 = x(2);
    GradDesc2_vals(i) = 200*(x2-x1^2)^2 + (1-x1)^2;
end

```

```

32     x = x0;
33     for i = 1:max_its
34         x = gd_backtrack(x, 1);
35         d_gd_backtrack(i) = norm(x - x_star);
36         x1 = x(1);
37         x2 = x(2);
38         gd_backtrack_vals(i) = 200*(x2-x1^2)^2 + (1-x1)^2;
39     end
40
41     it = 1:max_its;
42     figure(1);
43     plot(it, d_newton, 'b-', 'LineWidth', 2);
44     hold on;
45     plot(it, d_GradDesc2, 'r-', 'LineWidth', 2);
46     plot(it, d_gd_backtrack, 'g-', 'LineWidth', 2);
47     legend('Newton', 'Gradient Desc', 'GD Backtrack');
48     xlabel('iteration');
49     ylabel('||x - x*||');
50     title('distance to minimizer');
51     grid on;
52
53     figure(2);
54     semilogy(it, max(newton_vals, 1e-16), 'b-', 'LineWidth', 2);
55     hold on;
56     semilogy(it, GradDesc2_vals, 'r-', 'LineWidth', 2);
57     semilogy(it, gd_backtrack_vals, 'g-', 'LineWidth', 2);
58     legend('Newton', 'Gradient Desc', 'GD Backtrack');
59     xlabel('Iteration');
60     ylabel('f(x)');
61     title('value of f at each iteration');
62     grid on;

```



The plots show that Newton's method converges the quickest, converging in about 2 iterations. Then neither GD nor GD with backtracking converge to the minimizer  $[1,1]$  in 100 iterations, showing that both are slower to converge. GD with backtracking has a slightly faster than GD alone, but still slow, rate of convergence. GD alone has an extremely slow rate of convergence. Looking at the plot on the right, we can visualize the actual function values at each iteration as we supposedly get closer to the minimum, and the slopes for GD and GD with backtracking

3) we will redo the MNIST coding question but for differentiating 4s and 9s. Can reuse template for problem.

(a) Repeat parts (b), (c), and (e) from Q4 of HW4 for digits of 4s and 9s. Comment on the difference between the results and propose a reason as to why the performance did or did not change.

Command Window

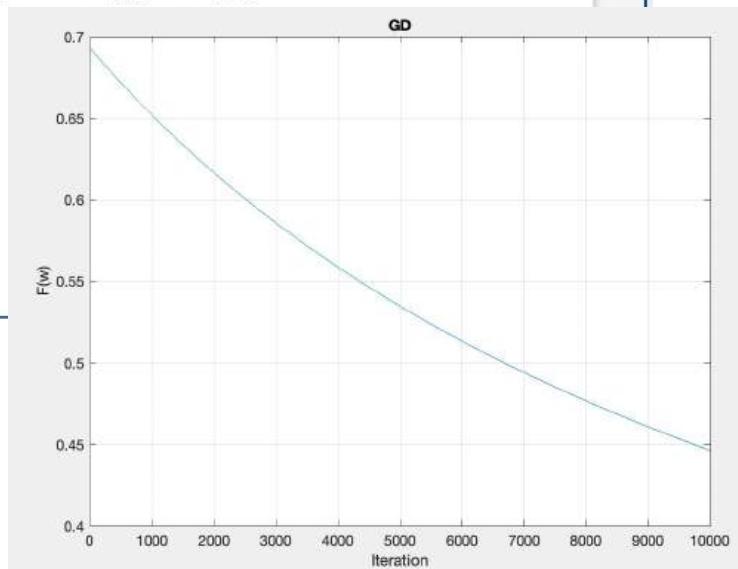
Editor - form\_data.m

```
1 form_data.m x mnist_train.csv x mnist_test.csv x untitled2 * +  
1 function [x,y] = form_data()  
2     data = readmatrix('mnist_train.csv');  
3     labels = table2array(data(:,1));  
4     pixels = table2array(data(:,2:end)) / 255.0;  
5  
6     fours_idx = find(labels == 4, 500);  
7     nines_idx = find(labels == 9, 500);  
8     x = [pixels(fours_idx,:); pixels(nines_idx,:)];  
9     y = [ones(500,1); -ones(500,1)];  
10    end
```

Command Window

Editor - GD\_hw5.m

```
+1 form_data.m x mnist_train.csv x mnist_test.csv x GD_hw5.m x GradDesc2.m x +  
1 function w = GD_hw5()  
2     [x,y] = form_data();  
3     [N, d] = size(x);  
4     u = 1e-4;  
5     max_its = 10000;  
6     F_vals = zeros(max_its,1);  
7     w = zeros(d,1);  
8  
9     for i = 1:max_its  
10        inner_prod = x*w;  
11        exp_term = exp(-y .* inner_prod);  
12        log_term = exp_term ./ (1 + exp_term);  
13        grad = -(1/N) * (x' * (y .* log_term));  
14        w = w - u * grad;  
15        F_vals(i) = (1/N) * sum(log(1 + exp_term));  
16    end  
17    figure;  
18    plot(1:max_its, F_vals);  
19    xlabel('Iteration');  
20    ylabel('F(w)');  
21    title('GD');  
22    grid on;  
23 end
```



```

24 %added to end of GD file from previous part
25 test_data = readtable('mnist_test.csv');
26 test_labels = table2array(test_data(:,1));
27 test_pixels = table2array(test_data(:,2:end)) / 255.0;
28 fours_idx_test = find(test_labels == 4, 500);
29 nines_idx_test = find(test_labels == 9, 500);
30 x_test = [test_pixels(fours_idx_test,:); test_pixels(nines_idx_test,:)];
31 y_test = [ones(500,1); -ones(500,1)];
32 classify_test = x_test * w;
33 test_pred = sign(classify_test);
34 test_err = sum(test_pred ~= y_test) / length(y_test);
35 train_pred = sign(x_train * w);
36 train_err = sum(train_pred ~= y_train) / length(y_train);
37
38 fprintf('Training error rate: %.4f\n', train_err);
39 fprintf('Test error rate: %.4f\n', test_err);
40 end

```

```

>> GD_hw5
Training error rate: 0.0810
Test error rate: 0.1260

```

Classification error rates shown above. As you can see, the error in the test data is slightly greater than the error in the training data, which is expected and usual, since the GD model was trained with the training data. The difference in errors for each data set is nothing alarming, its a small difference that is explained by the above. Both error percentages are relatively low. The performance in this run, when looking for distinctions between 4s and 9s, rather than 0s and 1s, is maybe not as good as the previous model. Meaning, the model is not as good at distinguishing between 4s and 9s as it is at distinguishing between 0s and 1s. However, this model had a more expected relationship between the training and testing errors, with the training error less than testing error, which the first model did not (possibly an issue with my implementation).

(b) Implement and run GD with backtracking line search. At each iteration  $t$ , initialize the step size  $\mu = 10^{-1}$ , and use  $\gamma = 0.5$  and  $\beta = 0.8$  to determine the correct  $\mu^{(t)}$ . Run your algorithm for at least 10000 iterations and plot the loss curve  $F(w^{(t)})$  as a function of  $t$ .

```

Editor - /Users/sabella/Downloads/GD_backtracking_hw5.m
+1 GD_hw5.m x mnist_train.csv x mnist_test.csv x GD_backtracking_hw5.m x +
1 function w = GD_backtracking_hw5()
2 [x,y] = form_data();
3 [N, d] = size(x);
4 max_its = 10000;
5 beta = 0.8;
6 gamma = 0.5;
7 F_vals = zeros(max_its,1);
8 w = zeros(d,1);
9
10 for i = 1:max_its
11     inner_prod = x * w;
12     exp_term = exp(-y .* inner_prod);
13     F_current = (1/N) * sum(log(1 + exp_term));
14     log_term = exp_term ./ (1 + exp_term);
15     grad = -(1/N) * (x' * (y .* log_term));
16     u = 0.1;
17     while true
18         w_new = w - u * grad;
19         inner_prod_new = x * w_new;
20         exp_term_new = exp(-y .* inner_prod_new);
21         F_new = (1/N) * sum(log(1 + exp_term_new));
22         if F_new <= F_current - gamma * u * (grad' * grad)
23             break;
24         end
25         u = beta * u;
26     end
27     w = w - u * grad;
28     F_vals(i) = F_current;
29 end
30 figure;
31 plot(1:max_its, F_vals, 'LineWidth', 2);
32 xlabel('Iteration');
33 ylabel('F(w^{(t)})');
34 title('GD Backtracking');
35 grid on;
36 end

```

∫

4) Again, we will read the MNIST coding question from HW4 but using different versions of GD  $w^{(t+1)} = w^{(t)} - \mu p^{(t)}$ . This time, we will differentiate 4s and 9s (we classify images of 4 and 9).

(a) Implement and run  $L^\infty$  GD with step size  $\mu = 10^{-8}$ . Run your alg for at least  $10^4$  iterations and initialize with  $w^{(0)} = 0$ . Recall:  $L^\infty$  GD uses steps  $p^{(t)} = \text{sign}(\nabla F(w^{(t)})) \|\nabla F(w^{(t)})\|_1$ .

```

Editor - /Users/sabella/Downloads/L_inf_GD.m
form_data.m mnist_train.csv mnist_test.csv L_inf_GD.m +
1 function w = L_inf_GD()
2     [x,y] = form_data();
3     [N, d] = size(x);
4     max_its = 10000;
5     u = 1e-8;
6     F_vals = zeros(max_its,1);
7     w = zeros(d,1);
8
9     for i = 1:max_its
10         inner_prod = x * w;
11         exp_term = exp(-y .* inner_prod);
12         log_term = exp_term ./ (1 + exp_term);
13         grad = -(1/N) * (x' * (y .* log_term));
14         grad_L1_norm = norm(grad, 1);
15         p = sign(grad) * grad_L1_norm;
16
17         w = w - u * p;
18         F_vals(i) = (1/N) * sum(log(1 + exp_term));
19     end
20     figure;
21     plot(1:max_its, F_vals, 'LineWidth', 2);
22     xlabel('Iteration');
23     ylabel('F(w^{(t)})');
24     title('L (infinity) GD');
25     grid on;
26 end
27

```

(b) Implement and run  $L^1$  GD with step size  $\mu = 10^{-4}$ . Run your alg for at least  $10^4$  iterations and initialize with  $w^{(0)} = 0$ . Recall:  $L^1$  GD uses steps  $p^{(t)} = \text{sign}(\nabla_j F(w^{(t)})) \|\nabla F(w^{(t)})\|_\infty e_{j^*}$ , where  $j^*$  is the location of the largest entry of the gradient, and  $e_j$  is the zero vector with a 1 in the  $j$ th entry.

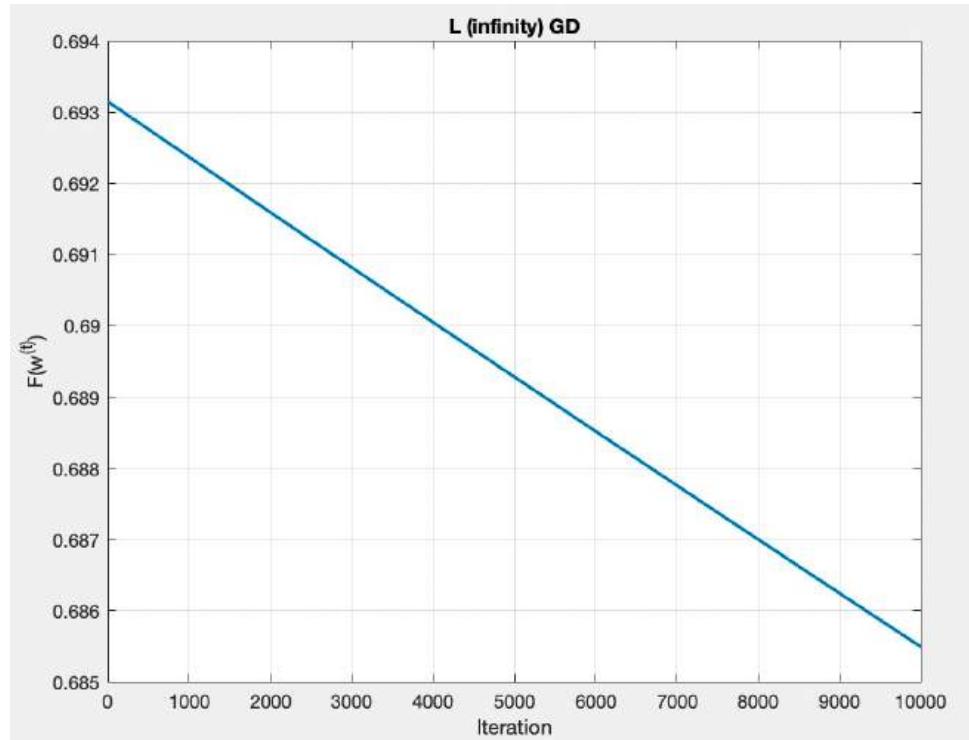
Editor - /Users/sabella/Downloads/L1\_GD.m \*

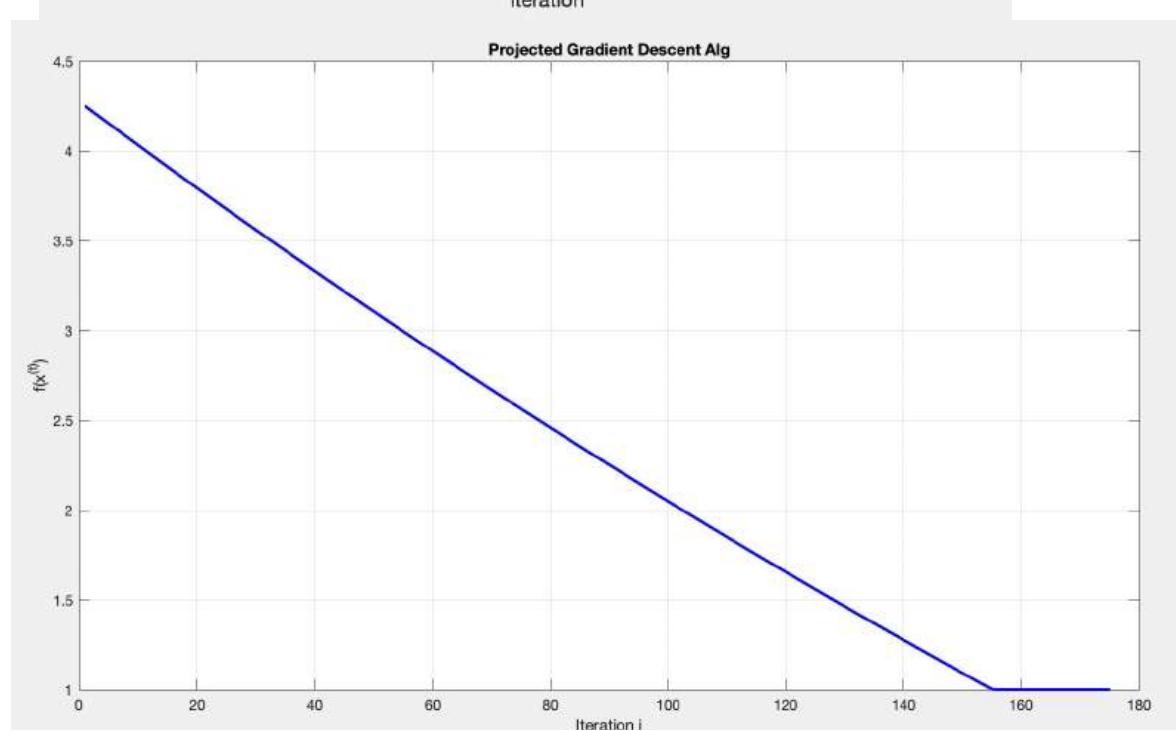
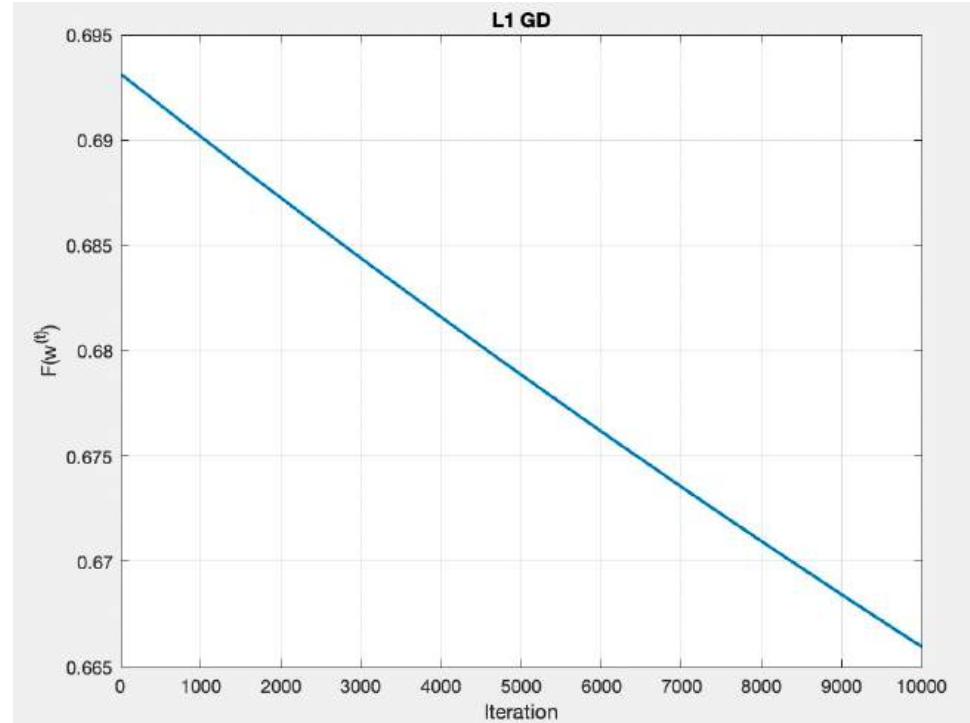
```

1 function w = L1_GD()
2     [x,y] = form_data();
3     [N, d] = size(x);
4     max_its = 10000;
5     mu = 1e-4;
6     F_vals = zeros(max_its,1);
7     w = zeros(d,1);
8
9     for i = 1:max_its
10        inner_prod = x * w;
11        exp_term = exp(-y .* inner_prod);
12        log_term = exp_term ./ (1 + exp_term);
13        grad = -(1/N) * (x' * (y .* log_term));
14
15        [~, j_star] = max(abs(grad));
16        grad_inf_norm = norm(grad, inf);
17        p = zeros(d,1);
18        p(j_star) = sign(grad(j_star)) * grad_inf_norm;
19
20        w = w - mu * p;
21        F_vals(i) = (1/N) * sum(log(1 + exp_term));
22    end
23    figure;
24    plot(1:max_its, F_vals, 'LineWidth', 2);
25    xlabel('Iteration');
26    ylabel('F(w^{(t)})');
27    title('L1 GD');
28    grid on;
29 end

```

(c) compare these two descent plots of  $F(w)$ , along with the analogous plot for GD from HW4. Which performs best, and do you have an argument for why? Do you think performance would change with different step sizes?





The GD implementation from HW4 appears to perform the best, and an argument for why could be that the regular GD always moves in the steepest direction of descent.  $L^\infty$  GD moves depending on each coordinate sign, which can possibly "confuse" the descent.  $L^2$  GD moves one coordinate at a time in the direction of  $j^*$ , which could be a problem since it didn't take into account all coordinates. Yes, I do think performance would change with different step sizes, because both  $L^\infty$  and  $L^2$  look to maybe have too small of a step and are converging too slowly.

All GDs rely on step size that essentially scale  $\hat{p}^{(t)}$  in the update. So, a larger step may lead to faster convergence for  $L^2$  &  $L^\infty$  GD. Clearly, the analogous GD converges faster by a large margin.

(a) For the coordinate descent problem, runn GD but store a running sum of which entry of  $p^{(t)}$  is nonzero at each iteration (not the actual value of the direction vector, just  $e_j$ ). This will result in a size 784 vector of mostly zeros, and should have integers at various entries whose sum equals the number of iterations. Reshape the vector to be a  $28 \times 28$  image and display the result. Why do you think these are the pixels that were chosen in the gradient? How can you use this to interpret the alg and its result?

```
+1 mnist_train.csv x mnist_test.csv x L_inf_GD.m x L1_GD.m x GradDesc.m x L1_GD_sum.m *
1 function coordinates = L1_GD_sum()
2 [x,y] = form_data();
3 [N, d] = size(x);
4 max_its = 10000;
5 u = 1e-4;
6 w = zeros(d,1);
7 coordinates = zeros(d,1);
8
9 for i = 1:max_its
10     inner_prod = x * w;
11     exp_term = exp(-y .* inner_prod);
12     log_term = exp_term ./ (1 + exp_term);
13     grad = -(1/N) * (x' * (y .* log_term));
14     [~, j_star] = max(abs(grad));
15     grad_inf_norm = norm(grad, inf);
16     p = zeros(d,1);
17     p(j_star) = sign(grad(j_star)) * grad_inf_norm;
18
19     w = w - u * p;
20     coordinates(j_star) = coordinates(j_star) + 1;
21 end
22 figure;
23 vector_image = reshape(coordinates, 28, 28);
24 imagesc(vector_image);
25 colorbar;
26 end
```

From the image, the words church were from rows ~15-17 and columns ~7-9. I think that these words was chosen due to the magnitude of their gradient being the largest, meaning that these words are the most influential in helping the model distinguish between  $y_s$  and  $q_s$ . This probably means that the algorithm repeatedly looks at this pixel in the data sets in order to help make decisions since it has the largest magnitude of gradient. Since our alg was L<sub>1</sub> GD, only one coord. updated at each iteration, so it makes sense that we just have one pixel for a small number of coordinates that get the large update from the large magnitude of gradient. So, the model sort of relies on a few coordinates and a small chunk of data to make its decisions, so it may not have the most accurate differentiation.

## Homework 4 for Math 173A - Fall 2025

For coding questions, you must submit your code and the requested answers. **Your code must be present to receive points.**

1. (a) Find an expression for the orthogonal projection of a point  $x \in \mathbb{R}^n$  onto the convex set

$$B = \{z \in \mathbb{R}^n : 0 \leq z_i \leq 1 \text{ for each } i = 1, \dots, n\}.$$

You need to show your work, and justify your answer. The expression can be written piecewise, and per dimension if it's easier / more compact.

**Hint:** It might be helpful to sketch  $B$ , when  $n = 2$  (i.e., in 2 dimensions), and use the sketch to help you figure out what the projection should be.

- (b) Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by

$$f(x) = \|Ax\|_2^2 + a^T x$$

where  $A \in \mathbb{R}^{n \times n}$  is a positive definite matrix, and  $a \in \mathbb{R}^n$ . Write a projected gradient descent algorithm to solve

$$\min_{x \in \Omega} f(x)$$

for  $\Omega = B$ , with  $B$  from part (a). You do not need to specify the step size for this problem.

- (c) Repeat part (b) but for  $\Omega = B_2^n = \{z \in \mathbb{R}^n : \|z\|_2 \leq 1\}$ .
2. Consider the *hollow* sphere  $S$  in  $\mathbb{R}^n$ , i.e., the set  $S := \{x \in \mathbb{R}^n : \|x\|_2^2 = 1\}$ . Consider the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$f(x) = x^T Q x$$

where  $Q$  is an  $n \times n$  symmetric matrix. For this problem you may use the fact that  $\nabla f(x) = 2Qx$ .

- (a) For an arbitrary point  $y \in \mathbb{R}^n$ ,  $\Pi(y)$  be the projection of  $y$  onto  $S$ . Find an expression for  $\Pi(y)$  and give a short argument (i.e., proof) for why this is the correct expression. Make sure to handle the case  $y = 0$  (i.e., the zero vector).

**Hint:** Recall that the projection minimizes  $\|x - y\|$  for  $y \in S$ . One approach would be to consider the reverse triangle inequality  $\|x - y\| \geq ||\|x\| - \|y\||$  and find a projection formula that achieves the global lower bound (i.e. equality instead of inequality). This would prove you've found the projection.

(b) Is  $S$  a convex set?

(c) Write a projected gradient descent algorithm, with constant step size  $\mu$ , for

$$\min_{x \in \mathbb{R}^n} x^T Q x \quad \text{subject to} \quad \|x\|_2^2 = 1.$$

(d) Is the projected gradient descent algorithm guaranteed to converge to the solution for small enough  $\mu$ ? If not, can you give an example of  $Q$  and an initialization  $x^{(0)}$  where the algorithm won't converge? **Hint:** Consider a diagonal matrix  $Q$  where not all entries are equal.

3. **Coding Question:** Consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  with  $f(x_1, x_2) = (x_1 - 2)^2 + (x_2 - 2)^2 - x_1 x_2$  and the following constrained optimization problem:

$$\min_{x_1, x_2} f(x_1, x_2) \quad \text{subject to} \quad 0 \leq x_i \leq 1, \quad i = 1, 2.$$

Program a projected gradient descent algorithm, with constant step size  $\mu = 0.001$  starting at  $(0.5, 0.5)$  for 175 iterations, for the above optimization problem. Submit your code and plot the function value  $f(x^{(t)})$  against the iteration  $t$ .

4. **Coding Question** For this question and the next question, you will need to download the MNIST data set. It is a data set of images of size 28x28 pixels. Each one is an image of a handwritten digit from 0-9. You can find some description of the dataset here: [https://en.wikipedia.org/wiki/MNIST\\_database](https://en.wikipedia.org/wiki/MNIST_database).

We've uploaded a Jupyter notebook template to the HW folder (HW4\_Q4.ipynb), and it provides a basic set of code for downloading / formatting MNIST and blocks for where to put your code / results / answers. You may need to twist the code a bit if it does not work for you.

We have also provided you directly the training and testing datasets, `mnist_train.csv` in the Assignment page. The `mnist_train.csv` file contains the 60,000 training examples and labels. The `mnist_test.csv` contains 10,000 test examples and labels. Each row consists of 785 values: the first value is the label (a number from 0 to 9) and the remaining 784 values are the pixel values (a number from 0 to 255).

In HW2 you considered the classification problem with logistic regression

$$F(w) = \frac{1}{N} \sum_{i=1}^N \log \left( 1 + e^{-\langle w, x_i \rangle y_i} \right). \quad (1)$$

You also wrote down a gradient descent algorithm for it.

- (a) Display one randomly selected image from your data (either the one from the template if you are using the template, or the training dataset `mnist_train.csv` if you directly use the dataset we provided) for each digit class. Provide the index number for each image.

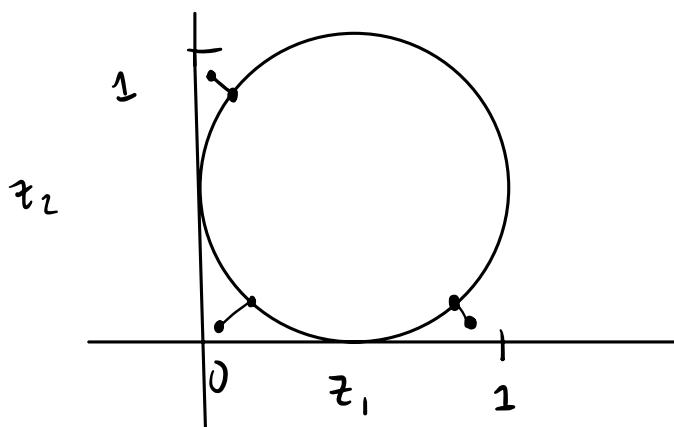
- (b) Continuing (a), select the first 500 examples of 0's and 1's for this example, those will form the training data  $(x_i, y_i) \in \mathbb{R}^{784} \times \{-1, 1\}$ ,  $i = 1, \dots, 1000$ . Assign label  $y_i = 1$  for 1s and  $y_i = -1$  for 0s. Also, renormalize your  $x_i$  so that the pixel values are floats between 0 and 1, instead of ints from 0 to 255.
- Note:** To get from images of size  $28 \times 28$  pixels to vectors in  $\mathbb{R}^{784}$ , you just need to “vectorize” the image. This means you can concatenate each of the 28 columns of the original image into one long vector of length 784. In Matlab, this is done with the command  $x(:)$ , similarly in Python it's `numpy.vectorize(x)` or `x.reshape()`.
- (c) Implement and run a Gradient Descent algorithm, with step-size  $\mu = 10^{-4}$ , to optimize the function (1) associated with this setup. You should run your algorithm for at least  $T = 10,000$  iterations, but if your computer can handle it, try  $T = 100,000$  or choose a reasonable stopping criterion. Provide a plot showing the value of  $F(w)$  at each iteration. Also, feel free to adjust  $\mu$  to be larger / smaller if the plot does not match your expectations.
- (d) Comment on the resulting plot. In particular, does the shape of  $F(w)$  suggest you've successfully converged to a local or global minimum? Does it appear you chose a good stopping criteria? Explain whether your answers to these questions are consistent with the theory we discussed in class (and in the notes). Be specific i.e., point to a specific theorem (or theorems) and indicate why it does or does not explain the behavior of the algorithm. Would the theory dictate a different choice of  $\mu$  than the one we used?
- (e) Now, use the  $w$  you found from part (c) to classify the next 500 data points associated to each of the 0 and 1 handwritten digits if you use the template we provided. The next 500 data points serve as test data. The idea is that you deploy your model on data you didn't use to train the model, as a check of its quality. If you use the datasets we provided, then you use the  $w$  you found from part (c) to classify the first 500 *test data points* associated to each of the 0 and 1 handwritten digits in `mnist.test.csv`. Recall that you need to use the function  $y = \text{sign}(w^T x)$  to classify. What was the classification error rate associated with the two digits on the test data (this should be a number between 0 and 1 )? What was it on the training data? Does this relationship make sense?

1) (a) Find an expression for the orthogonal projection of a point  $x \in \mathbb{R}^n$  onto the convex set

$$B = \{z \in \mathbb{R}^n : 0 \leq z_i \leq 1 \text{ for each } i=1,\dots,n\}$$

Show work, justify. Expression can be written piecewise and per dimension if it's easier/more compact Hint: sketch B, when  $n=2$  (2 dimensions), use sketch to figure out projection should be

$$\text{when } n=2 : B = \{z \in \mathbb{R}^2 : 0 \leq z_i \leq 1 \text{ for } i=1,2\}$$



$$\pi_B(x) \in \arg \min_{x \in B} \|x - z\|_2$$

$$\Rightarrow \pi_B(x) = \min_{\substack{x_i \in B \\ 0 \leq z_i \leq 1}} \sum_{i=1}^n (x_i - z_i)$$

$$z_i = \begin{cases} 0 & , x_i < 0 \\ x_i & , 0 \leq x_i \leq 1 \\ 1 & , x_i > 1 \end{cases}$$

Justification:  $(\pi_B(x) \in B)$

For any  $x$ ,  $0 \leq \pi_B(x) \leq 1$ . Thus,  $\pi_B(x) \in B$ .

Minimality:  $(\forall z \in B \quad |x - z| \geq \|\pi_B(x) - x\|_2)$

Suppose  $z \in \mathbb{R}^n : 0 \leq z_i \leq 1 \text{ for } i=1,\dots,n \quad (z \in B)$ .

Then RHS:  $\|x - z\|_2 = 0$

Suppose  $x < 0$ . Then LHS is abs val, so  $|x - z| \geq 0$  is true.

Suppose  $0 \leq x \leq 1$ . Then LHS is  $|x - z| = 0 \Rightarrow |0| \geq |0|$  is true.

Suppose  $x > 1$ . Then LHS is  $|x - z| \geq 0$  which is true since

$$(x > 1) - 1 \geq 0.$$

Set  $B$  is convex and closed, so  $\Pi_B(x)$  is unique.

(b) Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(x) = \|Ax\|_2^2 + a^T x$ , where  $A \in \mathbb{R}^{n \times n}$  is a positive definite matrix, and  $a \in \mathbb{R}^n$ . Write a projected GD algorithm to solve  $\min_{x \in \mathbb{R}^n} f(x)$  for  $\mathcal{S} = B$ , with  $B$  from part (a). You do not need to specify step size.

$$x^{(t+1)} = \arg \min_{x \in B} \|y^{(t+1)} - x\|_2, \quad y^{(t+1)} = x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})$$

$$\Rightarrow x^{(t+1)} = \Pi_B(x^{(t)} - \mu^{(t)} \nabla f(x^{(t)}))$$

$$f(x) = \|Ax\|_2^2 + a^T x = \sum_{i=1}^n (Ax)_i^2 + a^T x$$

$$\Rightarrow \nabla f(x) = \nabla \left( \sum_{i=1}^n (Ax)_i^2 + a^T x \right) = 2A^T A x + a$$

alg GD

Initialize  $x^{(0)} \in B$

(Initialize stepsize - not specified in problem)

for  $i = 0, 1, \dots$ , do

(update stepsize - not specified in this problem)

$$\text{grad\_it} = 2A^T A x^{(i)} + a$$

$$y^{(i+1)} = x^{(i)} - (\text{stepsize}) * \text{grad\_it}$$

If  $y^{(i+1)} < 0$ :

$$x^{(i+1)} = 0$$

else if  $y^{(i+1)} > 1$ :

$$x^{(i+1)} = 1$$

else:  $x^{(i+1)} = y^{(i+1)}$



$$(x^{(i+1)} = \Pi_B(y^{(i+1)}))$$

(c) Repeat part (b) but for  $N = B_2^n = \{z \in \mathbb{R}^n : \|z\|_2 \leq 1\}$

By a previous example,  $\Pi_{B_2^n}(x) = \begin{cases} x, & \|x\| \leq 1 \\ \frac{x}{\|x\|}, & \|x\| > 1 \end{cases}$

alg GD

Initialize  $x^{(0)} \in B$

(Initialize stepsize - not specified in problem)

for  $i = 0, 1, \dots, d_0$

(update stepsize - not specified in this problem)

$$\text{grad\_it} = 2A^T A x^{(i)} + a$$

$$y^{(i+1)} = x^{(i)} - (\text{stepsize}) * \text{grad\_it}$$

$$\text{mag\_y} = \|y^{(i+1)}\|_2 \quad // "2 norm of y^{(i+1)}$$

If  $\text{mag\_y} \leq 1$ :

$$x^{(i+1)} = y^{(i+1)}$$

$$\text{else: } x^{(i+1)} = y^{(i+1)} / \text{mag\_y}$$

$$\left. \begin{array}{l} x^{(i+1)} = \Pi_{B_2^n}(y^{(i+1)}) \end{array} \right\}$$

2) Consider the hollow sphere  $S$  in  $\mathbb{R}^n$ , i.e. the set  $S := \{x \in \mathbb{R}^n : \|x\|_2^2 = 1\}$ . Consider the function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  given by  $f(x) = x^T Q x$ , where  $Q$  is  $n \times n$  symmetric matrix. For this problem you may use the fact that  $\nabla f(x) = 2Qx$

(a) For an arbitrary point  $y \in \mathbb{R}^n$ ,  $\Pi(y)$  be the projection of  $y$  onto  $S$ . Find an expression for  $\Pi(y)$  and give a short argument (i.e. proof) for why this is the correct expression. Make sure to handle the case  $y = 0$  (i.e. the zero vector). Hint: Recall that the projection minimizes  $\|x - y\|$  for  $x \in S$ . One approach would be to consider the reverse triangle inequality  $\|x - y\| \geq \|x\| - \|y\|$  and find a projection formula that achieves the global lower bound (i.e. equality instead of inequality). This would prove you've found the projection.

$$\Pi_S(y) \in \arg \min_{x \in S} \|x - y\|_2$$

By reverse triangle inequality,  $\|x - y\|_2 \geq |\|x\| - \|y\||$ . We know that for any  $x \in S$ ,  $\|x\|_2^2 = 1 \Rightarrow \|x\|_2 = 1$

$$\Rightarrow \|x - y\|_2 \geq |1 - \|y\||$$

We want  $\|x - y\|_2 = |1 - \|y\||$  for any  $x \in S$ .

Starting from our goal  $|1 - \|y\||$ , we have

$$\begin{aligned} |1 - \|y\|| &= \left| \frac{\|y\|_2}{\|y\|_2} - \|y\|_2 \right| \\ &= \|y\|_2 \left| \left( \frac{1}{\|y\|_2} - 1 \right) \right| \\ &= \left\| \frac{1}{\|y\|_2} - y \right\|_2 \quad \Rightarrow \quad x = \frac{y}{\|y\|_2} \end{aligned}$$

$$\begin{aligned}
 \text{Projection: } \| \frac{y}{\|y\|_2} - y \|_2 &= \| y \left( \frac{1}{\|y\|_2} - 1 \right) \|_2 \\
 &= \|y\|_2 \left| \frac{1}{\|y\|_2} - 1 \right| \\
 &= |1 - \|y\|_2|
 \end{aligned}$$

Feasibility:  $(\Pi_S(y) \in S)$

$$\|x\|_2^2 = \left\| \frac{y}{\|y\|_2} \right\|_2^2 = 1 \quad \checkmark$$

When  $y = 0$ ,  $\|x - 0\|_2 \geq 1 \Rightarrow \|x\|_2 \geq 1$ , which holds true for any  $x \in S$ . Hence in this case any  $x \in S$  is a minimizer.

(b) Is  $S$  a convex set?

$S$  is the set of vectors in  $\mathbb{R}^n$  whose magnitude (length) squared equals one. Let  $u = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \in S$  and  $v = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \in S$  (both have magnitude 1). Then  $\alpha \begin{bmatrix} 1 \\ 0 \end{bmatrix} + (1-\alpha) \begin{bmatrix} 0 \\ -1 \end{bmatrix}$

$$= \begin{bmatrix} \alpha \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} \alpha \\ -1 \end{bmatrix}$$

Suppose  $\alpha = \frac{1}{2} \in [0, 1]$ . Then  $\alpha u + (1-\alpha)v = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$

$$\left\| \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \right\|_2 = \sqrt{\frac{1}{2}} \neq 1. \text{ So } \alpha u + (1-\alpha)v \notin S. \text{ Hence}$$

$S$  is not convex.

(c) Write a projected GD algorithm, with constant step size  $\mu$ , for  $\min_{x \in \mathbb{R}^n} x^T Q x$  subject to  $\|x\|_2^2 = 1$ .

We have  $\nabla f(x) = ZQx$ ,

$$x^{(t+1)} = \arg \min_{x \in S} \|y^{(t+1)} - x\|_2, \quad y^{(t+1)} = x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})$$

$$\Rightarrow x^{(t+1)} = \Pi_S(x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})) \text{ where}$$

alg GD

Initialize  $x^{(0)} \in S$

Initialize stepsize  $\nu$

for  $i = 0, 1, \dots$  do

$$\text{grad} = Z^T Q * x^{(i)}$$

$$y^{(i+1)} = x^{(i)} - \nu * \text{grad}$$

$$\text{mag-}y = \|y^{(i+1)}\|_2$$

If  $y^{(i+1)} = 0$ :

$$x^{(i+1)} = \text{any } x \in \mathbb{R}^n \text{ with } \|x\|_2 = 1$$

else:

$$x^{(i+1)} = y^{(i+1)} / \text{mag-}y$$

(d) Is the projected GD algorithm guaranteed to converge to the solution for small enough  $\mu$ ? If not, can you give an example of  $Q$  and an initialization  $x^{(0)}$  where the algorithm won't converge?  
Hint: consider a diagonal matrix  $Q$  where not all entries are equal.

Since  $S$  is not convex, the projected GD algorithm is not guaranteed to converge. Let's use

$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}, \quad x^{(0)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\Pi_S(y) = \begin{cases} \frac{y}{\|y\|_2}, & y \neq 0 \\ \text{any } x \in S, & y = 0 \end{cases}$$

$$\begin{aligned}
 \Rightarrow y^{(1)} &= x^{(0)} - 2\mu Q x^{(0)} \\
 &= \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 2\mu \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 2\mu & 0 \\ 0 & 4\mu \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 4\mu \end{bmatrix} = \begin{bmatrix} 0 \\ 1 - 4\mu \end{bmatrix}
 \end{aligned}$$

$$\|y^{(1)}\|_2 = \sqrt{(1-4\mu)^2} = |1-4\mu|$$

If  $\mu \neq 1/4$ ,  $y^{(0)} \neq 0$ , so

$$x^{(1)} = y^{(1)} / \|y^{(1)}\|_2 = \frac{\begin{bmatrix} 0 \\ 1 - 4\mu \end{bmatrix}}{|1 - 4\mu|} = \begin{bmatrix} 0 \\ \text{sign}(1 - 4\mu) \end{bmatrix}$$

Suppose  $0 < \mu < 1/4$ . Then  $1 - 4\mu > 0$  and thus

$$x^{(1)} = \begin{bmatrix} 0 \\ +1 \end{bmatrix} = x^{(0)}$$

So one iteration, even with  $0 < \mu < 1/4$ , does not move  $x^{(1)}$   
at all  $\Rightarrow$  does not converge

Suppose  $\mu > 1/4$ . Then  $1 - 4\mu < 0$ , so  $x^{(1)} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$

$$\begin{aligned}
 \Rightarrow y^{(2)} &= \begin{bmatrix} 0 \\ -1 \end{bmatrix} - \begin{bmatrix} 2\mu & 0 \\ 0 & 4\mu \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix} \\
 &= \begin{bmatrix} 0 \\ -1 \end{bmatrix} - \begin{bmatrix} 0 \\ -4\mu \end{bmatrix} = \begin{bmatrix} 0 \\ -(1 - 4\mu) \end{bmatrix}
 \end{aligned}$$

$$\Rightarrow x^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} = x^{(0)} \Rightarrow \text{does not converge}$$

3) Loading Question: Consider the function  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  with  
 $f(x_1, x_2) = (x_1 - 2)^2 + (x_2 - 2)^2 - x_1 x_2$  and the following constrained optimization problem:  $\min_{x_1, x_2} f(x_1, x_2)$  subject to  $0 \leq x_i < 1$ ,  $i=1,2$ .

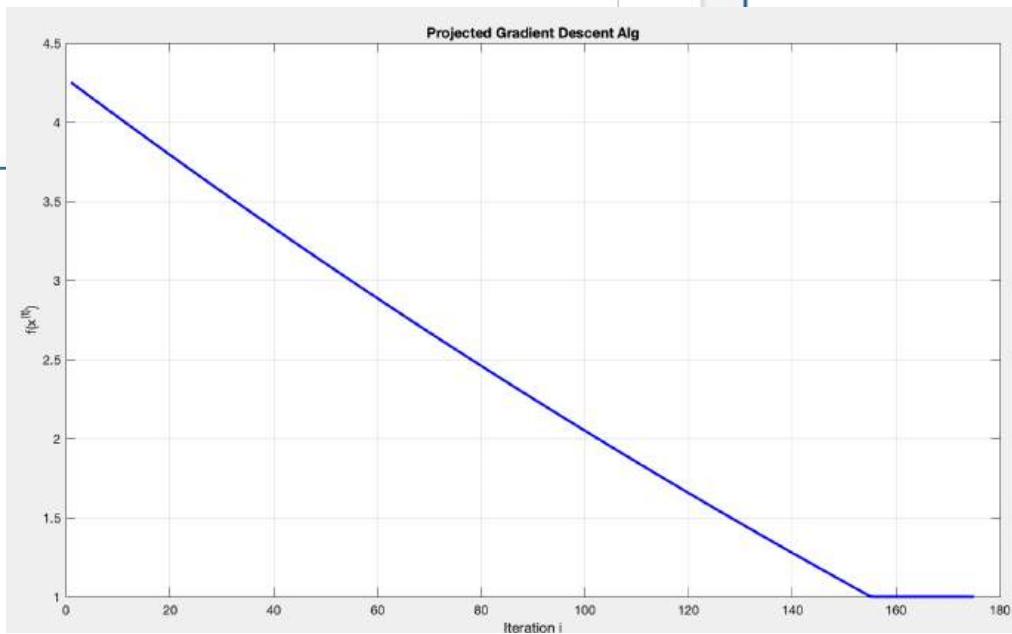
Program a projected GD alg, with constant step size  $\mu = 0.001$  starting at  $(0.5, 0.5)$  for 175 iterations, for the above. Submit code and plot function value  $f(x^{(t)})$  against iteration t.

```

Editor - /Users/sabella/Downloads/projected_GD.m
projected_GD.m + ✓

1 function projected_GD()
2     x = [0.5; 0.5];
3     u = 0.001;
4     max_it = 175;
5     f_values = zeros(max_it, 1);
6
7     for i = 1:max_it
8         f_values(i) = f(x);
9         g = grad(x);
10        x = x - u * g;
11        x = proj_f(x);
12    end
13
14    figure;
15    plot(1:max_it, f_values, 'b-', 'LineWidth', 2);
16    xlabel('Iteration i');
17    ylabel('f(x^{(t)})');
18    title('Projected Gradient Descent Alg');
19    grid on;
20
21 end
22
23 function val = f(x)
24     val = (x(1) - 2)^2 + (x(2) - 2)^2 - x(1)*x(2);
25 end
26
27 function g = grad(x)
28     g = [2*x(1) - x(2) - 4;
29             2*x(2) - x(1) - 4];
30 end
31
32 function x_proj = proj_f(x)
33     x_proj = max(0, min(1, x));
34 end

```

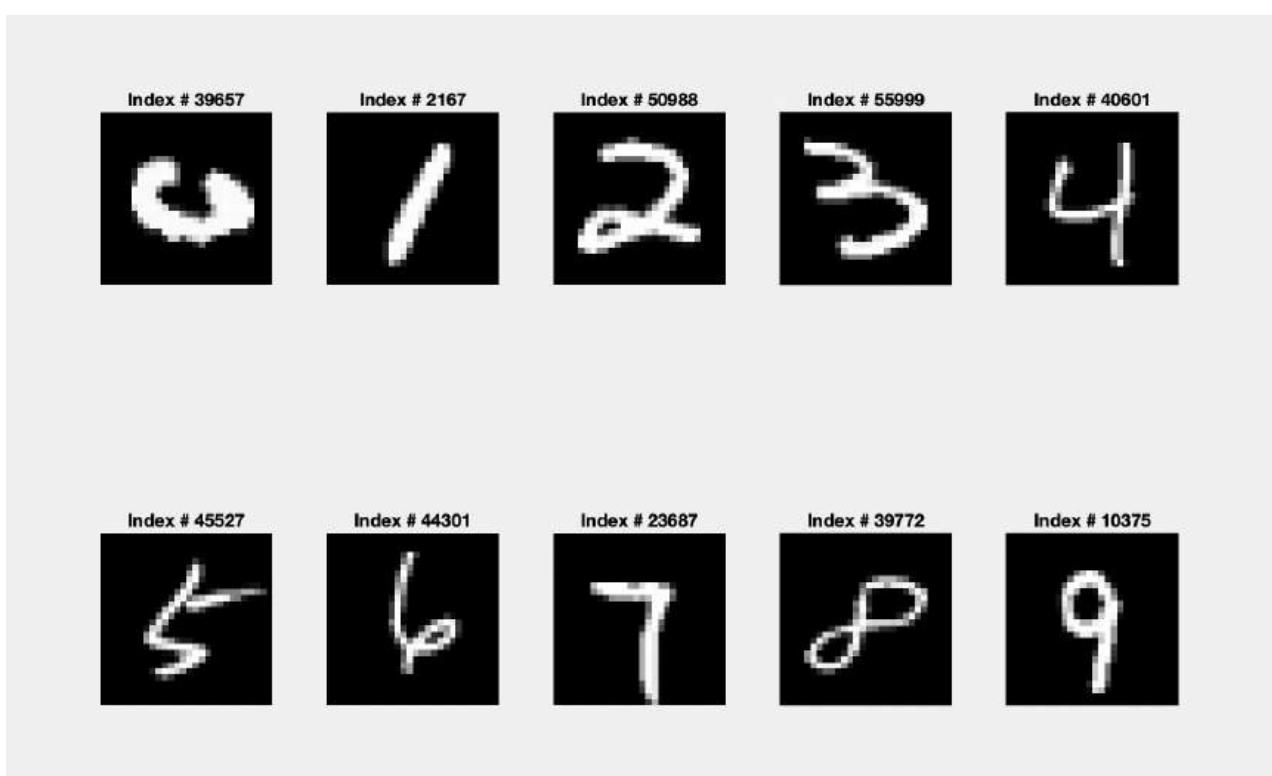


4) (a)

Editor - /Users/sabella/Downloads/display.m

projected\_GD.m display.m mnist\_train.csv +

```
1 function display()
2     data = readtable('mnist_train.csv');
3     labels = table2array(data(:,1));
4     pixels = table2array(data(:,2:end));
5
6     figure('Position', [100, 100, 1200, 600]);
7     for digit = 0:9
8         indices = find(labels == digit);
9         if isempty(indices), continue; end
10        random_idx = indices(randi(length(indices)));
11
12        img = reshape(pixels(random_idx, :), 28, 28)';
13        subplot(2, 5, digit+1);
14        imshow(img, []);
15        title(sprintf('Index # %d', random_idx));
16    end
17 end
```



(b)

Editor - /Users/sabella/Downloads/form.m

projected\_GD.m display.m mnist\_train.csv form.m +

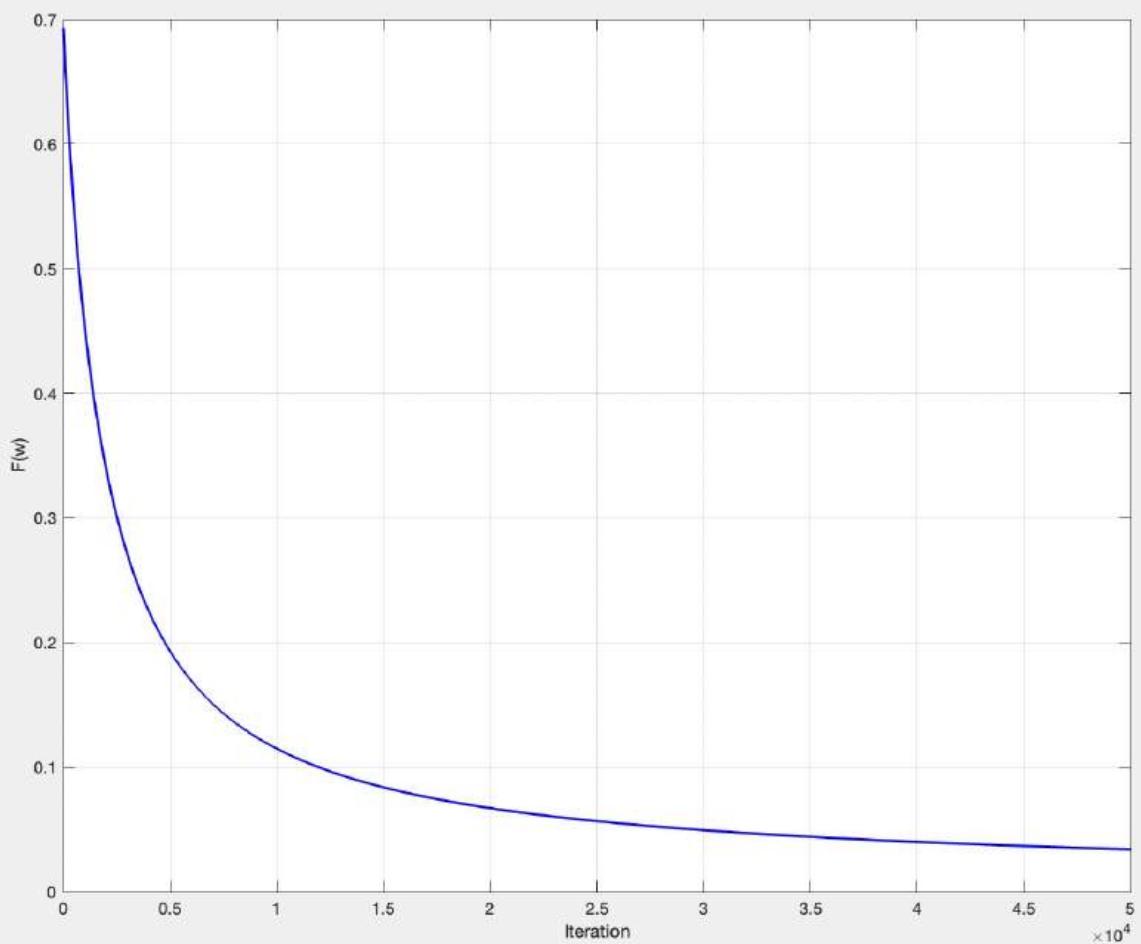
```
1 function [x, y] = form()
2     data = readtable('mnist_train.csv');
3     labels = table2array(data(:,1));
4     pixels = table2array(data(:,2:end)) / 255.0;
5
6     zeros_idx = find(labels == 0, 500);
7     ones_idx = find(labels == 1, 500);
8     x = [pixels=zeros_idx,:); pixels(ones_idx,:)];
9     y = [-ones(500,1); ones(500,1)];
10    end
```

Editor - /Users/sabella/Downloads/GradDesc.m

```

1 function [w, F_vals] = GradDesc()
2     [x, y] = form();
3     [~,d] = size(x);
4     u = 1e-4;
5     max_its = 50000;
6     F_vals = zeros(max_its, 1);
7
8     w = zeros(d, 1);
9
10    for i = 1:max_its
11        inner_prod = x * w;
12        F_vals(i) = mean(log(1 + exp(-y .* inner_prod)));
13        n = length(y);
14        denom = 1 + exp(y .* inner_prod);
15        grad = -(x' * (y ./ denom)) / n;
16
17        w = w - u * grad;
18    end
19
20    figure;
21    plot(1:max_its, F_vals, 'b-', 'LineWidth', 1.5);
22    xlabel('Iteration');
23    ylabel('F(w)');
24    grid on;
25 end

```



(d) The shape of  $F(w)$  does appear to show convergence to some minimum. At no point does the curve appear to increase at any iteration, and the curve descends steeply before rounding

out to a somewhat clear value. We found in HW2 that the function  $F$  is convex, and thus by a previous theorem any local minimizer is also a global minimizer. We also know that given a viable step size, GD converges to a global minimum, which explains the behavior of my GD alg. This theory does mention step size; I chose 50000 iterations since the step size of  $\mu = 10^{-9}$  seemed to not quite be enough to fully visualize the convergence. If we could show that  $F$  is Lipschitz (too much work for me right now) then a previous theorem would give us  $\mu \leq 1/L$ , which may lead to a quicker convergence and less iterations of GD. I believe I chose a good stopping criteria, since I knew 100000 would be a lot for my computer, but also 10000 did not seem to be enough. I chose an iteration count somewhere in the middle, and it visualizes the GD nicely.

(e)

```

Editor - /Users/sabella/Downloads/GradDesc.m
dish.m mnist_train.csv form.m GradDesc.m mnist_test.csv +
22
23     ylabel('1/w');
24     grid on;
25
26 %added to end of GD function from part (c)
27 classify = sign(x * w);
28 train_err = mean(classify ~= y);
29
30 data_test = readtable('mnist_test.csv'); %"form" but with test data
31 labels_test = table2array(data_test(:,1));
32 pixels_test = table2array(data_test(:,2:end)) / 255.0;
33 zeros_test = find(labels_test == 0, 500);
34 ones_test = find(labels_test == 1, 500);
35 x_test = [pixels_test=zeros_test,:]; pixels_test(ones_test,:);
36 y_test = [-ones(500,1); ones(500,1)];
37
38 classify_test = sign(x_test * w);
39 test_err = mean(classify_test ~= y_test);
40
41 fprintf('Training Error: %.4f, Test Error: %.4f\n', train_err, test_err);
42
43
44
Command Window
0

>> GradDesc
Training Error: 0.0030, Test Error: 0.0000

```

Classification error rates shown above. As you can see, the test error is 0, which is unusual and lower than the training error. I do find this relationship unusual, since I'd expect the test error to not fit the model as well as the training. However,

although unusual, this implies that the weight  $w^0$  applies to new data very, very well producing very no errors for this sample size. I think its possible that yes, the relationship is a bit unusual, but its due to the fact that this model is very accurate and trained well.

## Homework 3 for Math 173A - Fall 2025

1. Determine whether each function is Lipschitz, and if so find the smallest possible Lipschitz constant for the function. For all problems,  $\|\cdot\|$  represents the Euclidean norm (2-norm).
  - (a)  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \|x\|$
  - (b)  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \|x\|^2$
  - (c)  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  for  $\rho(x) = \frac{1}{1+e^{-x}}$ .
  - (d)  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \rho(w^T x + b)$  for some weight vector  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ , and  $\rho$  from part (c).
2. Let  $f$  be a convex and differentiable. Let  $x^*$  be the global minimum and suppose  $x^{(0)}$  is the initialization such that  $\|x^* - x^{(0)}\| \leq 5$ .
  - (a) Let  $f$  be  $L$ -Lipschitz function where  $L = 3$ . Determine the step size  $\mu$  and number of steps needed to satisfy
 
$$\left\| f\left(\frac{1}{t} \sum_{s=0}^{s-1} x^{(s)}\right) - f(x^*) \right\| \leq 10^{-4}.$$
  - (b) Let  $f$  be  $L$ -smooth where  $L = 3$ . Determine the step size  $\mu$  and number of steps needed to satisfy
 
$$\left\| f(x^{(t)}) - f(x^*) \right\| \leq 10^{-4}.$$
3. Consider the function  $f(x_1, x_2) = (2x_1 - 1)^4 + (x_1 + x_2 - 1)^2$ .
  - (a) Find the global minimum of  $f$ , and justify your answer.
  - (b) Starting at  $x^{(0)} = (0, 0)$ , perform gradient descent with backtracking line-search.
    - i. Starting at  $x^{(0)} = (0, 0)$  with stepsize, which is also called learning rate in the machine learning community,  $\mu^{(0)}$ , write down the gradient descent equation for  $x^{(1)}$ .
    - ii. Suppose we want to set  $\mu^{(0)}$  using backtracking line search with  $\gamma = 0.2$  and Armijo's condition  $f(x^{(1)}) \leq f(x^{(0)}) - \mu^{(0)}\gamma\|\nabla f(x^{(0)})\|_2^2$ . Find a value of  $\mu^{(0)}$  that satisfies this.
    - iii. Suppose instead you started with  $\mu^{(0)} = 1$  and an update of  $\mu^{(0)} \leftarrow \frac{1}{2}\mu^{(0)}$  (i.e.  $\beta = \frac{1}{2}$ ). In the worst case, how many steps of back-tracking would you have to take before accepting  $x^{(1)}$ ?

1) Determine whether each function is Lipschitz, if so find the smallest possible Lipschitz constant for the function. For all problems,  $\|\cdot\|$  represents the Euclidean norm (2-norm).

Lipschitz Condition:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\forall x, y \in \mathbb{R}^n$ ,  $L > 0$ :

$$|f(x) - f(y)| \leq L \cdot \|x - y\|_2$$

If  $f$  is  $C^2$ ,  $f$  is Lip iff  $\|\nabla f(x)\|_2 \leq L$

(a)  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \|x\|$

$$|f(x) - f(y)| = |\|x\|_2 - \|y\|_2| \leq \|x - y\|_2 \quad \forall x, y \in \mathbb{R}^n$$

by Reverse Triangle Inequality

$\Rightarrow$  Yes Lipschitz,  $L = 1$

$$\text{Wt } y=0, x \neq 0: \|f(x) - f(0)\| = \|\|x\|_2 - \|0\|_2\| = \|x\|_2$$

$$\Rightarrow \|x\|_2 \leq L \cdot \|x - 0\|_2 \Rightarrow \|x\|_2 \leq L \cdot \|x\|_2 \Rightarrow L \geq 1$$

$\Rightarrow L = 1$  smallest possible

(b)  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \|x\|^2$

$$|f(x) - f(y)| = |\|x\|_2^2 - \|y\|_2^2|$$

$$\|x\|_2^2 - \|y\|_2^2 = (x^T x) - (y^T y)$$

$$= (x + y)(x - y)$$

$$|(x + y)(x - y)| \leq \|x + y\| \|x - y\|$$

$$\Rightarrow |\|x\|_2^2 - \|y\|_2^2| \leq (\|x\|_2 + \|y\|_2) \|x - y\|_2$$

since  $x, y \in \mathbb{R}^n$ ,  $L = \|x\|_2 + \|y\|_2$  varies over  $\mathbb{R}^n$

$\Rightarrow f$  is NOT Lipschitz

$$(c) \rho: \mathbb{R} \rightarrow \mathbb{R} \text{ for } \rho(x) = \frac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$$

$$\begin{aligned}\nabla \rho(x) &= - (1+e^{-x})^{-2} (-e^{-x}) = \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \frac{e^{-x}}{(1+e^{-x})(1+e^{-x})} \\ &= \frac{1}{(1+e^{-x})} \left(1 - \left(\frac{1}{1+e^{-x}}\right)\right) \\ &= \rho(x)(1-\rho(x))\end{aligned}$$

$$\text{For } \rho \in (0,1) : f(\rho) = \rho(1-\rho) = \rho - \rho^2$$

$$\nabla f(\rho) = 1 - 2\rho = 0 \Rightarrow \rho = \frac{1}{2}$$

$$f\left(\frac{1}{2}\right) = \frac{1}{2}\left(1 - \frac{1}{2}\right) = \frac{1}{4}, \quad f(0) = 0, \quad f(1) = 0$$

$$\Rightarrow \max |\nabla \rho(x)| = \frac{1}{4} \Rightarrow L = \frac{1}{4} \quad \text{yes Lipschitz, smallest +}$$

(d)  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \rho(w^T x + b)$  for some weight vector  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ , and  $\rho$  from part (c).

$$|f(x) - f(y)| = |\rho(g(x)) - \rho(g(y))| \leq \frac{1}{4} \|\rho(g(x)) - \rho(g(y))\|_2$$

$$|g(x) - g(y)| = |(w^T x - b) - (w^T y - b)| = |w^T(x-y)|$$

by Cauchy-Schwarz  $|w^T(x-y)| \leq \|w\| \|x-y\|$

$$\Rightarrow L = \|w\|, \quad w \in \mathbb{R}^n$$

$$|f(x) - f(y)| \leq \frac{1}{4} \|w\| \|x-y\|_2$$

$$\Rightarrow L = \frac{1}{4} \|w\| \quad \text{yes Lipschitz, smallest +}$$

2. Let  $f$  be convex and differentiable. Let  $x^*$  be the global minimum and suppose  $x^{(0)}$  is the initialization such that  $\|x^* - x^{(0)}\| \leq 5$ .

(a) Let  $f$  be  $L$ -Lipschitz function where  $L=3$ . Determine the step size  $\mu$  and number of steps needed to satisfy

$$\left\| f\left(\frac{1}{t} \sum_{s=0}^{t-1} x^{(s)}\right) - f(x^*) \right\| \leq 10^{-4}$$

$$R=5, L=3, \|\nabla f\| \leq 3, f(y^{(t)}) - f(x^*) \leq \frac{RL}{\sqrt{t}}$$

$$\text{Choose } \mu = \frac{R}{L\sqrt{t}} = \frac{5}{3\sqrt{t}}$$

$$f(y^{(t)}) - f(x^*) \leq \frac{15}{\sqrt{t}} \quad \text{say} \leq 10^{-4}$$

$$\frac{15}{\sqrt{t}} \leq 10^{-4} \Rightarrow \frac{15}{10^{-4}} \leq \sqrt{t} \quad \text{Let } \sqrt{t} = 150000 \\ \Rightarrow t = \underline{2.25 \times 10^{10}}$$

$$\Rightarrow \mu = \frac{5}{3 \cdot 150000} \approx \underline{1.11 \times 10^{-5}}$$

(b) Let  $f$  be  $L$ -smooth where  $L=3$ . Determine the step size  $\mu$  and number of steps needed to satisfy

$$\|f(x^{(t)}) - f(x^*)\| \leq 10^{-4}.$$

$$\underline{\mu \in (0, \frac{1}{3})} \quad f(x^{(t)}) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\mu t} = 10^{-4}, \forall t \geq 1$$

$$\Rightarrow \|x_0 - x^*\|^2 = 2 \cdot 10^{-4} \mu t \quad \Rightarrow \mu = \frac{\|x_0 - x^*\|^2}{2 \cdot 10^{-4} t}$$

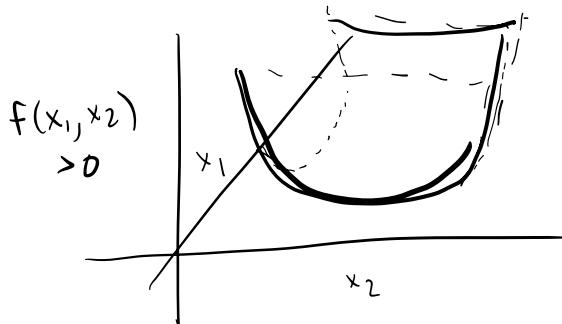
$$\text{Let } \boxed{\mu = \frac{1}{3}} = \frac{\|x_0 - x^*\|^2}{2 \cdot 10^{-4} t} \quad \Rightarrow t = \frac{3 \cdot \|x_0 - x^*\|^2}{2 \cdot 10^{-4}}, \|x_0 - x^*\|^2 \leq 25$$

$$\Rightarrow t \geq \frac{75}{2 \cdot 10^{-4}} \Rightarrow t \geq 375000$$

$$\mu = \frac{1}{3}, t = 375000$$

3. Consider the function  $f(x_1, x_2) = (2x_1 - 1)^4 + (x_1 + x_2 - 1)^2$

(a) Find the global min of  $f$ , justify your answer



$$f(x_1, x_2) \geq 0 \quad \forall x_1, x_2$$

$\Rightarrow$  find min by setting  
 $f(x_1, x_2) = 0$

$$(2x_1 - 1)^4 = 0 \Rightarrow 2x_1 - 1 = 0 \Rightarrow x_1 = \frac{1}{2}$$

$$\Rightarrow 2\left(\frac{1}{2} + x_2 - 1\right) = 0 \Rightarrow 2\left(x_2 - \frac{1}{2}\right) = 0 \Rightarrow 2x_2 - 1 = 0 \Rightarrow x_2 = \frac{1}{2}$$

$$\Rightarrow x^* = \left(\frac{1}{2}, \frac{1}{2}\right)$$

$$\nabla(2x-1)^4 = 8(2x-1)^3 \quad \nabla^2(2x-1)^4 = 48(2x-1)^2 > 0 \Rightarrow \text{convex}$$

$$\nabla(x_1 + x_2 - 1)^2 = \begin{bmatrix} 2(x_1 + x_2 - 1) \\ 2(x_1 + x_2 - 1) \end{bmatrix} \quad \nabla^2(x_1 + x_2 - 1)^2 = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} = \text{PD} \Rightarrow \text{convex}$$

(Prin. Theorem:  $x^*$  is global min iff  $\nabla f(x^*) = 0$ )

$$\nabla f = \begin{pmatrix} 8(2x-1)^3 + 2(x_1 + x_2 - 1) \\ 2(x_1 + x_2 - 1) \end{pmatrix}$$

$$\nabla f\left(\frac{1}{2}, \frac{1}{2}\right) = \begin{pmatrix} 8(0)^3 + 2\left(\frac{1}{2} + \frac{1}{2} - 1\right) \\ 2\left(\frac{1}{2} + \frac{1}{2} - 1\right) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \checkmark$$

$\Rightarrow x^* = \left(\frac{1}{2}, \frac{1}{2}\right)$  is a global min since  $f(x_1, x_2)$  is the sum of two non-negative convex functions and  $\nabla f\left(\frac{1}{2}, \frac{1}{2}\right) = (0, 0)$

(b) Starting at  $x^{(0)} = (0, 0)$ , perform gradient descent with backtracking line-search.

i. Starting at  $x^{(0)} = (0, 0)$  with stepsize, aka learning rate in ML community,  $\mu^{(0)}$ , write down the gradient descent equation for  $x^{(1)}$ .

$$\nabla f(x_1, x_2) = \begin{bmatrix} 8(2x_1 - 1)^3 + 2(x_1 + x_2 - 1) \\ 2(x_1 + x_2 - 1) \end{bmatrix}$$

$$\nabla f(0, 0) = \begin{bmatrix} 8(-1)^3 + 2(-1) \\ 2(-1) \end{bmatrix} = \begin{bmatrix} -10 \\ -2 \end{bmatrix}$$

$$\begin{aligned} x^{(1)} &= x^{(0)} - \mu^{(0)} \nabla f(x^{(0)}) \\ &= (0, 0) - \mu^{(0)} \cdot \begin{bmatrix} -10 \\ -2 \end{bmatrix} = \mu^{(0)} \begin{bmatrix} 10 \\ 2 \end{bmatrix} = \begin{bmatrix} 10\mu^{(0)} \\ 2\mu^{(0)} \end{bmatrix} \end{aligned}$$

ii. Suppose we want to set  $\mu^{(0)}$  using backtracking line search with  $\gamma = 0.2$  and Armijo's condition

$$f(x^{(1)}) \leq f(x^{(0)}) - \mu^{(0)} \gamma \|\nabla f(x^{(0)})\|_2^2.$$

Find a value of  $\mu^{(0)}$  that satisfies this.

$$f(x^{(0)}) = f(0, 0) = (2(0) - 1)^4 + (0 + 0 - 1)^2 = 1 + 1 = 2$$

$$\|\nabla f(x^{(0)})\|_2^2 = \left\| \begin{bmatrix} -10 \\ -2 \end{bmatrix} \right\|_2^2 = 100 + 4 = 104$$

$$\begin{aligned} \Rightarrow f(x^{(1)}) &\leq 2 - \mu^{(0)} \cdot 0.2 \cdot 104 \\ &\leq 2 - 20.8\mu^{(0)} \end{aligned}$$

$$\Rightarrow f(\mu^{(0)} \begin{bmatrix} 10 \\ 2 \end{bmatrix}) \leq 2 - 20.8\mu^{(0)}$$

$$\begin{aligned} f(\mu^{(0)} \begin{bmatrix} 10 \\ 2 \end{bmatrix}) &= f\left(\begin{bmatrix} 10\mu^{(0)} \\ 2\mu^{(0)} \end{bmatrix}\right) = (2(10\mu^{(0)}) - 1)^4 + (10\mu^{(0)} + 2\mu^{(0)} - 1)^2 \\ &= (20\mu^{(0)} - 1)^4 + (12\mu^{(0)} - 1)^2 \end{aligned}$$

$$\Rightarrow (20\mu^{(0)} - 1)^4 + (12\mu^{(0)} - 1)^2 \leq 2 - 20.8\mu^{(0)}$$

$$\text{Test } \mu^{(0)} = 1: (20-1)^4 + (12-1)^2 \leq 2 - 20.8 \\ \Rightarrow (19)^4 + 11^2 \leq 2 - 20.8 \Rightarrow \text{F}$$

Test  $\mu^{(0)} = 0.01$ :

$$(20 \cdot 0.01 - 1)^4 + (12(0.01) - 1)^2 \leq 2 - 20.8(0.01)$$

$$\Rightarrow (0.2 - 1)^4 + (0.12 - 1)^2 \leq 2 - 0.208$$

$$\Rightarrow (-.8)^4 + (-.88)^2 \leq 1.792$$

$$\Rightarrow 0.4096 + 0.7744 \leq 1.792$$

$$\Rightarrow 1.184 \leq 1.792 \quad \checkmark \quad \mu^{(0)} = 0.01 \text{ satisfies}$$

iii. Suppose instead you started with  $\mu^{(0)} = 1$  and an update of  $\mu^{(0)} \leftarrow \frac{1}{2} \mu^{(0)}$  (ie.  $\beta = \frac{1}{2}$ ). In the worst case, how many steps of back-tracking would you have to take before accepting  $x^{(1)}$ ?

$\mu^{(0)}$  cuts in half each iteration, i.e.

$$\mu^{(0)} = \frac{1}{2}, \mu^{(1)} = \frac{1}{4}, \mu^{(2)} = \frac{1}{8}, \dots \\ \Rightarrow \mu^{(t)} = \frac{1}{2^t}$$

Have  $(20\mu^{(0)} - 1)^4 + (12\mu^{(0)} - 1)^2 \leq 2 - 20.8\mu^{(0)}$  satisfied

$$\text{by } \mu^{(0)} = 0.01 \Rightarrow \frac{1}{2^t} \leq 0.01$$

$$\Rightarrow 2^t \geq 100$$

$$\Rightarrow t \geq 7$$

$$2^7 = 128 \Rightarrow \mu^{(t)} = \frac{1}{128} = 0.0078125$$

$$\Rightarrow t = 7$$

## Homework 2 for Math 173A - Fall 2025

1. Using the conditions of optimality, find the local maximizers and local minimizers, of the following functions and determine whether they are local maximizers or local minimizers. You may use a computer to find the eigenvalues, but these questions should have easily accessible eigenvalues by hand.

- (a)  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  for  $f(x_1, x_2) = x_1^4 + 2x_2^4 - 4x_1x_2$   
 (b)  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  for  $f(\vec{x}) = \vec{x}^T A \vec{x} + b^T \vec{x}$ , where

$$A = \begin{bmatrix} -1 & 0 & 1/2 \\ 0 & -1 & 0 \\ 1/2 & 0 & -1 \end{bmatrix}, \quad b = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

Bonuses question: Are the local maximizers/minimizers you find also global maximizers/minimizers? If so, why?

2. Consider the problem  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \|Ax - b\|_2^2$  for  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ . (Note: this was on the last homework). Write down the gradient descent algorithm to solve the optimization

$$\min_{x \in \mathbb{R}^n} f(x).$$

This doesn't have to be a computer program, just something of the form

$$\begin{aligned} x^{(0)} &= \dots \\ x^{(t+1)} &= \dots \text{ (where the right hand side is in terms of } x^{(t)}\text{).} \end{aligned}$$

You need to specify the formula of the gradient, but you do not need to specify the exact choice of the stepsize.

3. **Implementing Classification Model:** First some background for classification:

- You are given labeled data  $\{(x_i, y_i)\}_{i=1}^N$  for  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$ .
- Logistic regression involves choosing a label according to

$$y = \text{sign}(\langle w, x \rangle).$$

Note we ignore the y-intercept term here, so we only need the optimal  $w \in \mathbb{R}^d$ .

- It turns out the correct function to minimize to find the weights is

$$F(w) = \frac{1}{N} \sum_{i=1}^N \log \left( 1 + e^{-\langle w, x_i \rangle y_i} \right).$$

**Questions:**

- Is  $F(w)$  a convex function?
- Write down the gradient descent algorithm for minimizing  $F$ . You need to specify the starting point, the formula of the gradient, but you do not need to specify the exact choice of the stepsize.

4. **Coding Question:** Recall that the equation for an ellipse in  $\mathbb{R}^2$  is

$$a_1 x^2 + a_2 y^2 = 1.$$

Given data  $\{(x_i, y_i)\}_{i=1}^N \subset \mathbb{R}^2$  that lie on (or near) the ellipse, you can find the best fit ellipse by solving

$$\min_{a \in \mathbb{R}^2} f(a)$$

where

$$f(a) = \sum_{i=1}^N (a_1 x_i^2 + a_2 y_i^2 - 1)^2$$

- Find an  $A \in \mathbb{R}^{N \times 2}$  and  $b \in \mathbb{R}^N$  such that  $f(a) = \|Aa - b\|_2^2$ . What is  $A$  in terms of  $(x_i, y_i)$ ?
- Download the data provided on the HW page (called HW2\_ellipse.csv), and create a scatter plot of the points (submit your code and the plot).
- Using Problem 2, create computer code to compute the gradient descent algorithm on this  $f(a)$ . The code must include a stopping condition. Use a step-size of  $\eta = \frac{1}{2\|A^T A\|}$ . Note, you cannot use a built-in gradient descent algorithm it must be written with a while or for loop. Also note, the norm of a matrix  $\|X\| = \lambda_{\max}(X)$  is the largest eigenvalue (in magnitude) of  $X$ , and can be computed using “norm(X,2)” in matlab or “np.linalg.norm(X,2)” in Python. (submit the code)
- Using the data provided and your gradient descent code, estimate the solution  $a$ . Report  $a$  and  $f(a)$ . Given  $f(a)$  and  $N$ , do you think you fit the data well or poorly? Given the convexity of  $f$ , do you think this is the optimal  $a$ ?

1) Using the conditions of optimality, find the local maximizers and minimizers, of the following function & determine whether they are local maximizers or local minimizers.

$$(a) f: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ for } f(x_1, x_2) = x_1^4 + 2x_2^4 - 4x_1x_2$$

$$\nabla f(x) = 0 \Rightarrow \nabla f(x_1, x_2) = \begin{pmatrix} 4x_1^3 - 4x_2 \\ 8x_2^3 - 4x_1 \end{pmatrix} = 0$$

$$4x_1^3 - 4x_2 = 0 \Rightarrow 4x_1^3 = 4x_2 \Rightarrow x_1^3 = x_2$$

$$8x_2 - 4x_1 = 0 \Rightarrow 8x_2^3 - 4x_1$$

$$\text{using } x_1^3 = x_2 : 8(x_1^3)^3 - 4x_1 = 8x_1^9 - 4x_1 = 0$$

$$= 2x_1^9 - x_1 = 0$$

$$= x_1(2x_1^8 - 1) = 0$$

$$\Rightarrow x_1 = 0 \Rightarrow x_2 = 0, \quad 2x_1^8 = 1 \Rightarrow x_1^8 = \frac{1}{2} \Rightarrow x_1 = \pm \left(\frac{1}{2}\right)^{1/8}$$

$$x_1^* = (0, 0)$$

$$x_2^* = \left( \left(\frac{1}{2}\right)^{1/8}, \left(\frac{1}{2}\right)^{3/8} \right)$$

$$x_3^* = \left( -\left(\frac{1}{2}\right)^{1/8}, -\left(\frac{1}{2}\right)^{3/8} \right)$$

$$\nabla^2 f(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} \end{bmatrix} = \begin{bmatrix} 12x_1^2 & -4 \\ -4 & 24x_2^2 \end{bmatrix}$$

$$\nabla^2 f(0, 0) = \begin{bmatrix} 0 & -4 \\ -4 & 0 \end{bmatrix} \begin{bmatrix} -\lambda & -4 \\ -4 & -\lambda \end{bmatrix} = \lambda^2 - 16 = 0$$

$$\Rightarrow \lambda_1 = -4, \lambda_2 = +4$$

$\Rightarrow$  saddle pt

$$\nabla^2 f\left(\left(\frac{1}{2}\right)^{1/8}, \left(\frac{1}{2}\right)^{3/8}\right) = \begin{bmatrix} 12\left(\frac{1}{2}\right)^{1/4} & -4 \\ -4 & 24\left(\frac{1}{2}\right)^{3/4} \end{bmatrix} \quad (12\left(\frac{1}{2}\right)^{1/4} > 0) \checkmark$$

$$\det(\nabla^2 f) = (12\left(\frac{1}{2}\right)^{1/4} \cdot 24\left(\frac{1}{2}\right)^{3/4}) - 16 \\ = 288 \cdot \left(\frac{1}{2}\right) - 16 = 144 - 16 = 128 > 0$$

$$\nabla^2 f\left(-\left(\frac{1}{2}\right)^{1/8}, -\left(\frac{1}{2}\right)^{3/8}\right) = \begin{bmatrix} 12\left(\frac{1}{2}\right)^{1/4} & -4 \\ -4 & 24\left(\frac{1}{2}\right)^{3/4} \end{bmatrix} \quad \text{same as above}$$

$\nabla^2 f\left(\left(\frac{1}{2}\right)^{1/8}, \left(\frac{1}{2}\right)^{3/8}\right)$  &  $\nabla^2 f\left(-\left(\frac{1}{2}\right)^{1/8}, -\left(\frac{1}{2}\right)^{3/8}\right)$  both PD,  
 so  $x_2^* = \left(\left(\frac{1}{2}\right)^{1/8}, \left(\frac{1}{2}\right)^{3/8}\right)$  &  $x_3^* = \left(-\left(\frac{1}{2}\right)^{1/8}, -\left(\frac{1}{2}\right)^{3/8}\right)$   
 are local minima

(b)  $f: \mathbb{R}^3 \rightarrow \mathbb{R}$  for  $f(\vec{x}) = \vec{x}^T A \vec{x} + b^T \vec{x}$ , where

$$A = \begin{bmatrix} -1 & 0 & 1/2 \\ 0 & -1 & 0 \\ 1/2 & 0 & -1 \end{bmatrix}, \quad b = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \quad f(\vec{x}) = \frac{1}{2} \vec{x}^T (2A) \vec{x} + b^T \vec{x} \\ \Rightarrow \nabla f(\vec{x}) = 2A\vec{x} + b = 0 \\ \Rightarrow 2Ax = -b \quad \Rightarrow Ax = -\frac{1}{2}b = -\frac{1}{2} \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ -1/2 \\ -1 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} -1 & 0 & 1/2 \\ 0 & -1 & 0 \\ 1/2 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ -1/2 \\ -1 \end{bmatrix}$$

$$\Rightarrow -x_1 + \frac{1}{2}x_3 = 0 \quad \Rightarrow x_1 = \frac{1}{2}x_3 \\ -x_2 = -\frac{1}{2} \quad \Rightarrow x_2 = \frac{1}{2}$$

$$\frac{1}{2}x_1 - x_3 = -1$$

$$\Rightarrow \frac{1}{2} \left( \frac{1}{2} x_3 \right) - x_3 = -1$$

$$\Rightarrow -\frac{3}{4} x_3 = -1 \Rightarrow x_3 = \frac{4}{3} \Rightarrow x_1 = \frac{2}{3}$$

$$\Rightarrow x^* = \begin{bmatrix} 2/3 \\ 1/2 \\ 4/3 \end{bmatrix}$$

$$\nabla^2 f(\vec{x}) = 2A = \begin{bmatrix} -2 & 0 & 1 \\ 0 & -2 & 0 \\ 1 & 0 & -2 \end{bmatrix}$$

$$\det \begin{bmatrix} -2-\lambda & 0 & 1 \\ 0 & -2-\lambda & 0 \\ 1 & 0 & -2-\lambda \end{bmatrix} = (-2-\lambda)((-2-\lambda)^2 - 1) = 0$$

$$\Rightarrow -2-\lambda = 0 \Rightarrow \underline{\lambda_1 = -2}$$

$$(-2-\lambda)^2 = 1 \Rightarrow -2-\lambda = \pm 1$$

$$+1 : -2-\lambda = 1 \Rightarrow \underline{\lambda_2 = -3}$$

$$-1 : -2-\lambda = -1 \Rightarrow \underline{\lambda_3 = -1}$$

$\lambda_1, \lambda_2, \lambda_3 < 0$  so  $\nabla^2 f$  neg det  $\Rightarrow$

$$x^* = \begin{bmatrix} 2/3 \\ 1/2 \\ 4/3 \end{bmatrix} \text{ is local maximizer}$$

(Prin. Theorem:  $x^*$  is global min iff  $\nabla f(x^*) = 0$ )

$$(a) x_2^* = \left( \left(\frac{1}{2}\right)^{1/8}, \left(\frac{1}{2}\right)^{3/8} \right) \quad x_3^* = \left( -\left(\frac{1}{2}\right)^{1/8}, -\left(\frac{1}{2}\right)^{3/8} \right)$$

$$\nabla f(x_2^*) = \begin{pmatrix} 4\left(\frac{1}{2}\right)^{3/8} - 4\left(\frac{1}{2}\right)^{3/8} \\ 8\left(\frac{1}{2}\right)^{9/8} - 4\left(\frac{1}{2}\right)^{1/8} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$x_2^* = \left( \left(\frac{1}{2}\right)^{1/8}, \left(\frac{1}{2}\right)^{3/8} \right)$  is a global min

$$\nabla f(x_2^*) = \begin{pmatrix} -4\left(\frac{1}{2}\right)^{3/8} + 4\left(\frac{1}{2}\right)^{3/8} \\ -8\left(\frac{1}{2}\right)^{9/8} + 4\left(\frac{1}{2}\right)^{1/8} \end{pmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$x_3^* = \left( -\left(\frac{1}{2}\right)^{1/8}, -\left(\frac{1}{2}\right)^{3/8} \right)$  is a global min

(ignore)

2) Consider the problem  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  for  $f(x) = \|Ax - b\|_2^2$  for  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ . Write down the gradient descent algorithm and solve the optimization  $\min_{x \in \mathbb{R}^n} f(x)$ .

Doesn't have to be comp program, just something of form  
 $x^{(0)} = \dots$

$$x^{(t+1)} = \dots \text{ (terms of } x^{(t)})$$

Specify formula of gradient, do not need to spec exact choice of step size.

$$\text{From last hw: } \nabla f(x) = 2A^T(Ax - b)$$

### Gradient Descent (GD)

Choose  $x^{(0)} \in \mathbb{R}^n$

For  $t = 1, 2, \dots$ , (or until stopping criterion met)

$$\text{set } x^{(t)} = x^{(t-1)} - \mu^{(t-1)} \nabla f(x^{(t-1)})$$

$$(\text{Based on Piazza Q21}) \quad x^{(t)} = x^{(t-1)} - \mu^{(t-1)} \nabla f(x^{(t-1)}) \quad \text{for } t \geq 0$$

Let  $x^{(0)} = 0$ . (chosen initialization)

$$\begin{aligned} \Rightarrow x^{(1)} &= x^{(0)} - \mu^{(0)} \nabla f(x^{(0)}) \\ &= 0 - \mu^{(0)} (2A^T(A \cdot 0 - b)) \\ &= 0 + \underline{\mu^{(0)} \cdot 2A^T b} = \underline{2\mu A^T b} \end{aligned}$$

$$\begin{aligned} x^{(2)} &= x^{(1)} - \mu^{(1)} \nabla f(x^{(1)}) \\ &= \mu^{(0)} 2A^T b - \mu^{(1)} (2A^T (A(\mu^{(0)} 2A^T b) - b)) \\ &= \mu^{(0)} 2A^T b - \mu^{(1)} (2A^T (\mu^{(0)} 2AA^T b - b)) \\ &= \cancel{\mu^{(0)} 2A^T b} - \mu^{(1)} (\cancel{4\mu^{(0)} A^T A A^T b} - \cancel{2A^T b}) \\ &= \underline{4\mu A^T b - 4\mu^2 (A^T A) A^T b} \end{aligned}$$

⋮

$$x^{(t+1)} = x^{(t)} - \mu (2A^T (A x^{(t)} - b))$$

3) You are given labeled data  $\{(x_i, y_i)\}_{i=1}^N$  for  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$ .

Logistic regression involves choosing a label according to

$$y = \text{sign}(\langle w, x \rangle).$$

Note we ignore y-intercept term here, only need optimal  $w \in \mathbb{R}^d$ .

Correct function to minimize to find the weights is

$$F(w) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{-\langle w, x_i \rangle y_i}).$$

(a) Is  $F(w)$  a convex function?

$$\nabla f = \frac{\partial}{\partial \langle w, x_i \rangle y_i} \log(1 + e^{-\langle w, x_i \rangle y_i})$$

$$= \frac{1}{1 + e^{-\langle w, x_i \rangle y_i}} (e^{-\langle w, x_i \rangle y_i})' = \frac{-e^{-\langle w, x_i \rangle y_i}}{1 + e^{-\langle w, x_i \rangle y_i}} = \frac{-1}{1 + e^{\langle w, x_i \rangle y_i}}$$

$$\nabla^2 f = \frac{\partial^2}{\partial \langle w, x_i \rangle y_i} \left( \frac{-1}{1 + e^{\langle w, x_i \rangle y_i}} \right) = \frac{e^{\langle w, x_i \rangle y_i}}{(1 + e^{\langle w, x_i \rangle y_i})^2} > 0$$

$\forall w, x_i, y_i \in \mathbb{R}$

$\sin u e^x > 0 \quad \forall x \in \mathbb{R}$ .

$\sin u \nabla^2 f > 0$ ,  $f$  is convex. The sum of convex functions is

also convex, so  $F = \frac{1}{N} \sum_{i=1}^N f$  is convex.

(b) Write down gradient descent alg for minimizing  $F$ . You must specify starting point, formula of gradient, not step size  $\mu$ .

Choose  $w^{(0)} \in \mathbb{R}^n$   
 For  $t = 1, 2, \dots$ , (or until stopping criterion met)  
 set  $w^{(t+1)} = w^{(t)} - \mu^{(t)} \nabla F(w^{(t)})$

$$\nabla F(w) = \frac{1}{N} \sum_{i=1}^N \nabla (\log(1 + e^{-\langle w, x_i \rangle y_i}))$$

$$\begin{aligned} & \text{(ignore norm)} \quad \frac{1}{1 + e^{-\langle w, x_i \rangle y_i}} \cdot e^{-\langle w, x_i \rangle y_i} \cdot (-x_i y_i) \\ & = \end{aligned}$$

$$= \frac{-x_i y_i}{1 + e^{\langle w, x_i \rangle} y_i} \Rightarrow \nabla F(w) = \frac{1}{N} \sum_{i=1}^N \left( \frac{-x_i y_i}{1 + e^{\langle w, x_i \rangle} y_i} \right)$$

## Gradient Descent (GD)

$w^{(0)} = 0$  // choose  $w^{(0)} \in \mathbb{R}^n$

For  $t = 1, 2, \dots$ ,

$$w^{(t+1)} = w^{(t)} - \mu^{(t)} \cdot \frac{1}{N} \sum_{i=1}^N \left( \frac{-x_i y_i}{1 + e^{\langle w^{(t)}, x_i \rangle} y_i} \right)$$

4) Coding: Recall that the equation for ellipse in  $\mathbb{R}^2$  is  
 $a_1 x^2 + a_2 y^2 = 1$ .

Given data  $\{(x_i, y_i)\}_{i=1}^N \subset \mathbb{R}^2$  that lie on (or near) the ellipse, you can find the best fit ellipse by solving  $\min_{a \in \mathbb{R}^2} f(a)$ , where

$$f(a) = \sum_{i=1}^N (a_1 x_i^2 + a_2 y_i^2 - 1)^2$$

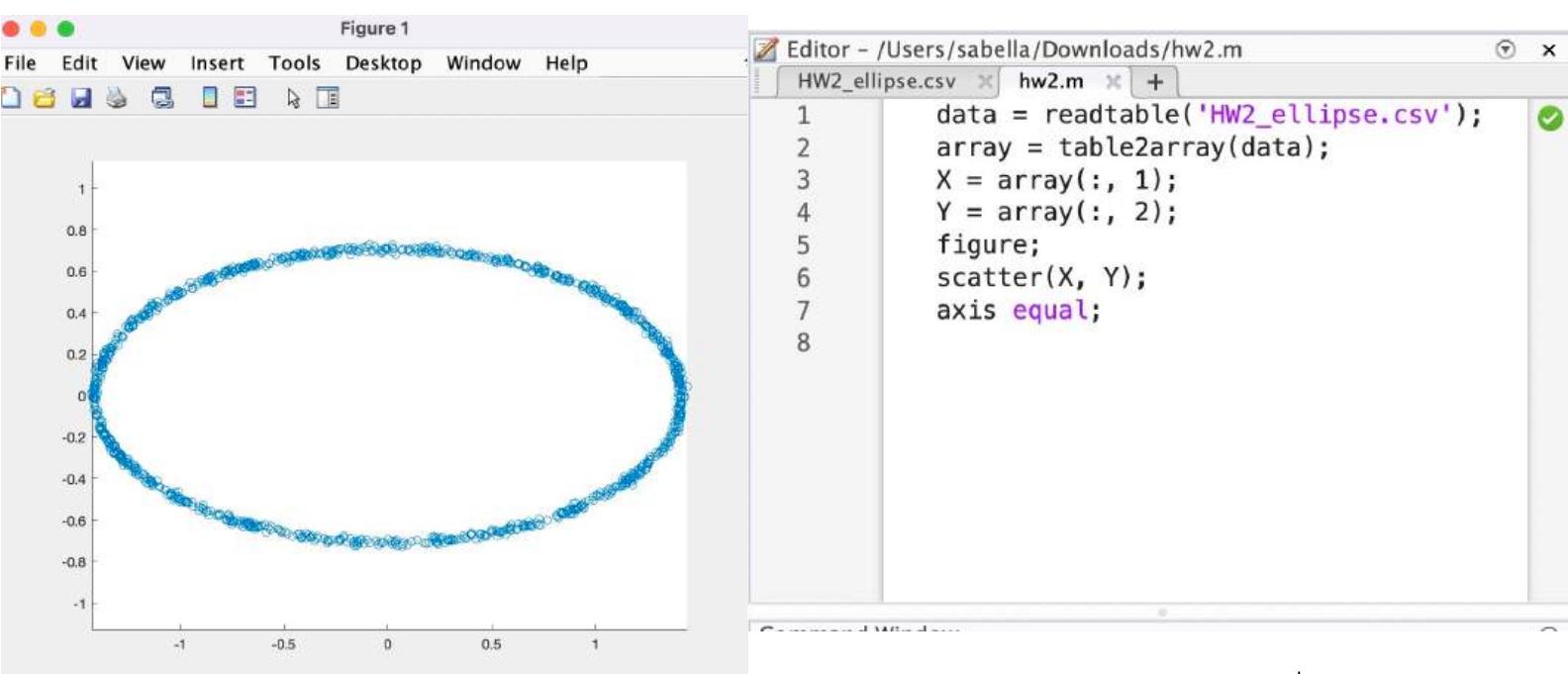
(a) Find an  $A \in \mathbb{R}^{N \times 2}$  and  $b \in \mathbb{R}^N$  such that  $f(a) = \|Aa - b\|_2^2$ . What is  $A$  in terms of  $(x_i, y_i)$ ?

$$\begin{aligned} f(a) &= \sum_{i=1}^N (a_1 x_i^2 + a_2 y_i^2 - 1)^2 = \|a_1 x_i^2 + a_2 y_i^2 - 1\|_2^2 \\ &= \left\| \underbrace{\begin{bmatrix} x_1^2 & y_1^2 \\ x_2^2 & y_2^2 \\ \vdots & \vdots \\ x_N^2 & y_N^2 \end{bmatrix}}_{A \in \mathbb{R}^{N \times 2}} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} - 1 \right\|_2^2 \end{aligned}$$

for  $i = 1, \dots, N$  :

$$\begin{aligned} &\left[ \begin{array}{cc} x_1^2 & y_1^2 \\ x_2^2 & y_2^2 \\ \vdots & \vdots \\ x_N^2 & y_N^2 \end{array} \right] \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \\ &\Rightarrow A (\in \mathbb{R}^{N \times 2}) \quad a \quad - \quad b (\in \mathbb{R}^N) \end{aligned}$$

(b) Download the data provided, create scatter plot & submit code & plot



(c) Using Problem 2, create computer code to compute the gradient descent alg on this  $f(a)$ . The code must include a stopping condition. Use a step-size of  $\eta = \frac{1}{2\|A^T A\|}$ . Cannot use built-in GD alg, must involve for/while loop. Norm of a matrix  $\|X\| = \lambda_{\max}(X)$ , can be computed using "norm(X,2)" in matlab

```

Editor - /Users/sabella/Downloads/hw2.m
HW2_ellipse.csv  hw2.m  +
1 data = readtable('HW2_ellipse.csv');
2 array = table2array(data);
3 X = array(:, 1);
4 Y = array(:, 2);
5 N = length(X);
6
7 A = [X.^2, Y.^2];
8 b = ones(N,1);
9
10 a = zeros(2,1);
11 max_its = 1000;
12 tol = 1e-6;
13 step = 1 / (2 * norm(A'*A, 2));
14
15 for t = 1:max_its
16     del_f = 2 * A' * (A*a - b);
17     a_t = a - step * del_f;
18
19     if norm(a_t - a) < tol
20         break;
21     end
22
23     a = a_t;
24 end
25 disp(a);
26
27 figure;
28 scatter(Y, Y).

```

Name
a
A
a_t
array
b
data
del_f
max_its
N
step
t
tol
X
Y

(d) Using the data provided and your gradient descent code, estimate the solution  $a$ . Report  $a$  and  $f(a)$ . Given  $f(a)$  and  $N$ , do you think you fit the data well or poorly? Given the convexity of  $f$ , do you think this is the optimal  $a$ ?

The estimate for  $a$  I found given the data and my code is

$$a = \begin{bmatrix} 1/2 \\ 2 \end{bmatrix}.$$

Command Window

```
>> hw2
[0.5001
 1.9946] = a
0.4641 = f(a)
fx >>
```

$$\Rightarrow f(a) = \sum_{i=1}^N (a_1 x_i^2 + a_2 y_i^2 - 1)^2$$

$$f\left(\begin{bmatrix} 0.5001 \\ 1.9946 \end{bmatrix}\right) = \sum_{i=1}^{1000} \underbrace{(0.5001 x_i^2 + 1.9946 y_i^2 - 1)^2}_{\text{we want this}} = 0.4641$$

$$= \sum_{i=1}^{1000} ((\text{estimate}) - (\text{actual}))^2$$

$$= \sum_{i=1}^{1000} (\text{total error})^2 = 0.4641$$

$$\Rightarrow \frac{f(a)}{N} = \frac{0.4641}{1000} = \frac{\text{total squared error}}{\# \text{ points}} = \underline{0.0004641}$$

The average squared error per point is 0.0004641, which is considerably small compared to the values of  $(x_i^2, y_i^2)$  from observing the CSV file. Hence I think I fit the data well.

Since  $f(a) = \|Aa - b\|_2^2$ , we found previously (HW 1, problem 5c and problem 6a) that  $f(a)$  is convex. Hence, there is only one global minimizer that minimizes  $f(a)$ . We found this  $a$  by using GD on  $f(a)$ , so yes, I do believe this  $a$  is optimal.