

HW07- CS 189

1.

Who else did you work with on this homework? In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

I worked on this homework with Ehimare Okoyomon, Prashanth Ganeth, and Daniel Mockaitis. We worked by getting together throughout the week and communicating on facebook.

I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up}

Nicholas Lorio, 26089160

Question 6 Own Question.

Where did machine learning neural networks get their name?

Neural networks are loosely inspired by the biological neural networks of animal brains. These systems progressively learn and improve their functions for certain tasks based off of observational data and feedback. This is opposed to task specific programming.

Our brains take in observational data and process it abstractly. Certain neurons fire, detecting abstract features of the world around us. Its a grouping of abstract features/patterns.

Example Observation:

Yuor abiilty to exaimne hgiher-lveel ffeaures is waht aollws yuo to unedrtsand tihs.

Neural networks in machine learning use affine tools to interpret and build models from observational data / inputs .

HW2 - 189.

2.a. $x, b \in \mathbb{R}^d$. Prove $f(x) = \|x - b\|_2$ is a convex function of x

$$f(x) = \sqrt{(x-b)^T(x-b)} = \sqrt{(x^T - b^T)(x-b)}$$

$$\nabla f(x) = \frac{1}{2} \left(x^T x - x^T b - b^T x + b^T b \right)^{-1/2} \left(2x - b - b \right)$$

$$\begin{aligned} \nabla_x(b^T x), \quad b^T(x+\Delta) &= b^T x + b^T \Delta & b^T(b^T x)/\delta x &= b^T \quad \nabla_x(b^T x) = (b^T)^T = b \\ \nabla_x(x^T x) & \quad (x+\Delta)^T(x+\Delta) = (x^T + \Delta^T)(x+\Delta) = x^T x + x^T \Delta + \Delta^T x + \Delta^T \Delta \\ &= (x^T x) \Delta + (\Delta^T x)^T + (\Delta^T \Delta)^T \\ &= (x^T x - x^T b - b^T x + b^T b)^{-1/2} (x - b) \end{aligned}$$

$$\nabla_x^2 = \begin{cases} -((x-b)^T(x-b))^{-3/2} & (x-b)^T(x-b) \\ +2((x-b)^T(x-b))^{-1/2} & \end{cases}$$

• If $\nabla_x^2(f(x)) \geq 0$ then our function $f(x)$ is convex.

$$\begin{aligned} \nabla_x^2 &= -((x-b)^T(x-b))^{(-3/2)} (x-b)^T(x-b) + ((x-b)^T(x-b))^{-1/2} \\ &= -((x-b)^T(x-b))^{-1/2} + 2((x-b)^T(x-b))^{-1/2} \end{aligned}$$

$$= -\frac{\|x-b\|_2^2}{\|x-b\|_2^3} + \frac{1}{\|x-b\|_2} \geq 0 \geq 0 \quad \text{as required for convexity}$$

2.b. Gradient Descent $x_0 = [0, 0]$ $x \in \mathbb{R}^2$, $b \in [4.5, 6]$, $t_i = 1$
 $\epsilon = 0.01$, want $\|x_i - x^*\|_2 \leq \epsilon$

$$x_{i+1} = x_i - \nabla f(x_i) \quad \text{w/ } t_i = 1$$

• Gradient descent will work for our convex function.

$$QF(x) = ((x-b)^T(x-b))^{-1/2} (x-b)^T = ((x - [4.5])^T(x - [4.5]))^{-1/2} (x - [4.5])$$

• Set $QF(x)$ equal to zero to solve for our optimal solution.
 Convex means this local min is a global min.

$$\frac{1}{\|x-b\|_2} (x-b)^T = 0 \quad \frac{x}{\|x-b\|_2^2} = \frac{b}{\|x-b\|_2^2} \quad x^* = b = [4.5] \\ [6]$$

$$\bullet x_1 = [6] - QF(0) = [6] + \left(\left(\begin{bmatrix} 4.5 \\ 6 \end{bmatrix}^T \cdot \begin{bmatrix} 4.5 \\ 6 \end{bmatrix} \right)^{-1/2} \begin{bmatrix} -4.5 \\ -6 \end{bmatrix} \right) \stackrel{-1/2}{=} (-56.25) \begin{bmatrix} 4.5 \\ 6 \end{bmatrix}$$

$$= \begin{bmatrix} +0.6 \\ +0.8 \end{bmatrix}$$

$$\bullet x_2 = \begin{bmatrix} +0.6 \\ +0.8 \end{bmatrix} \quad \nabla x F(x_1) =$$

$$2.c. \quad \|x_{i+1}\|_2 = \|x_i - \alpha F(x_i)\|_2 = \|x_i - (5/6)^i F(x_i)\|_2$$

$$x_2 = x_0 + 1 \cdot OF(x_0) = 0 + \frac{b}{\|b\|}, \quad OF(x_0) = \frac{b}{\|b\|}$$

$$x_2 = x_1 - (5/6) OF(x_1) = x_1 + \frac{5/6}{\|b\|} \frac{b}{\|b\|} = \frac{b}{\|b\|} (1 + 5/6)$$

$$x_3 = \frac{b}{\|b\|} \left(\sum_{i=0}^{A-1} (5/6)^i \right) = \frac{b}{\|b\|} \cdot \frac{1}{1 - 5/6} = \frac{b}{\|b\|} \cdot 6$$

$$\left[x_3 = \frac{b}{\|b\|} \cdot 6 \right]$$

- converges to a point.
- if optimal point is further away than our converge point than this step size for bisection method will never reach the optimal point.

enough to our starting No.

$$\text{2.d. } x_1 = x_0 - \frac{1}{2} \nabla F(x_0), \quad \nabla F(x_0) = ((0-b)^T(0-b))^{-1/2}(0-b) = \frac{-b}{\|b\|_2} = \frac{-b}{\|b\|}$$
$$= 0 - \nabla F(x_0) = +b/\|b\|$$

$$x_2 = x_1 - \frac{1}{2} \nabla F(x_1), \quad \nabla F(x_1) = \left(\left(\frac{b}{\|b\|} - b \right)^T \left(\frac{b}{\|b\|} - b \right) \right)^{-1/2} \left(\frac{b}{\|b\|} - b \right)$$
$$= b \left(\frac{1}{\|b\|} - 1 \right) / \|b \left(\frac{1}{\|b\|} - 1 \right)\|_2$$
$$= b \left(1 - \frac{1}{\|b\|} \right) / \|b \left(1 - \frac{1}{\|b\|} \right)\|_2 = \frac{-b}{\|b\|}$$

$$x_2 = \frac{b}{\|b\|} + \frac{1}{2} \frac{b}{\|b\|} = \frac{b}{\|b\|} \left(1 + \frac{1}{2} \right)$$

$$x_1 = \frac{b}{\|b\|} \left(\sum_{i=1}^k \frac{1}{i} \right)$$

2.e. cont.

$$x_i = b/\|b\| \left(\sum_{j=1}^i \frac{1}{g_j} \right)$$

$$\|x_i - x^*\|_2 \leq 0.01, \left\| \frac{b}{\|b\|} \sum_{j=1}^i \frac{1}{g_j} - b \right\| \leq 0.01$$

$$\left\| \frac{b(\log(i) - 1)}{\|b\|} \right\|_2 \leq 0.01, \left(\frac{\log(i) - 1}{\|b\|} \right) \|b\| \leq 0.01$$

$$\log(i) - \|b\| \leq 0.01, \log(i) \leq 0.01 + \|b\|$$

$$i \geq \exp(0.01 + \|b\|) \quad | \quad \begin{array}{l} \text{\# iterations for convergence to} \\ \text{optimal solution} \end{array}$$

3.

$$\|U^T V\| \leq \|U\|_2 \|V\|_2$$

A.T does not

$$F(x) = \frac{1}{2} \|Ax - b\|_2^2$$

$$a. b = \vec{0}$$

Step size: γ , each $x \in \mathbb{R}^d$ state

$$F(x) = \frac{1}{2} \|Ax - b\|_2^2 = \frac{1}{2} (Ax - b)^T (Ax - b) = \frac{1}{2} (x^T A^T Ax - x^T A^T b - b^T A x + b^T b) = \frac{1}{2} x^T A^T A x - b^T A x + \frac{1}{2} b^T b$$

$$\nabla_x F(x) = \frac{1}{2} ((A^T A + A^T A)x - 2A^T b) = A^T A x - A^T b$$

$$x_{i+1} = x_i - \gamma \nabla F(x_i), \quad x_1 = x_0 - \gamma \nabla F(x_0) = x_0 - \gamma (A^T A x_0 - A^T b) \\ = (I - \gamma A^T A) x_0$$

$$x_2 = x_1 - \gamma \nabla F(x_1) = x_0 - \gamma A^T A x_0 + \gamma A^T b - \gamma \nabla F(x_1)$$

$$= x_0 - \gamma (A^T A x_0 - A^T b) - \gamma (A^T A (x_0 - \gamma (A^T A x_0 - A^T b)) - A^T b)$$

$$= x_0 - \gamma A^T A x_0 - \gamma A^T A y_1 + \gamma A^T A y_1$$

$$= x_0 - \gamma (2A^T A x_0 - A^T A y_1 A^T A x_0) = (I - \gamma 2A^T A + \gamma A^T A y_1 A^T A) x_0 \\ = (I - \gamma A^T A)^2 x_0$$

$$x_3 = x_2 - \gamma (A^T A (x_0 - \gamma (2A^T A x_0 - A^T A y_1 A^T A x_0)))$$

$$= x_0 - \gamma (2A^T A x_0 - A^T A y_1 A^T A y_1) - \gamma (A^T A (x_0 - \gamma (2A^T A x_0 - A^T A y_1 A^T A x_0)))$$

$$x_n = (I - \gamma A^T A)^n x_0$$

3.b. all eig $B \leq I$ in absolute value.

$$B = I - \gamma A^T A \text{ from 3.a.}$$

• $A^T A$ is a PSD matrix so $\lambda_i > 0$ ①

• γ is a scalar > 0 ②

① $\therefore A^T A = P D P^T, \quad I = P P^T P P^T$

thus $I - \gamma A^T A = P (P^T P - \gamma D) P^T$

$$= P \begin{pmatrix} 1 - \gamma \lambda_0 \\ \vdots \\ 1 - \gamma \lambda_n \end{pmatrix} P^T$$

• we need $\|1 - \gamma \lambda_i\| \leq 1 \quad \forall i \quad \left\{ \begin{array}{l} 1 - \gamma \lambda_i \geq -1 \\ 1 - \gamma \lambda_i \leq 1 \end{array} \right. \quad \left. \begin{array}{l} -\gamma \lambda_i \geq -1 \\ -\gamma \lambda_i \leq 1 \end{array} \right\} -\gamma \lambda_i \geq -1$

• $\lambda_i > 0, \gamma > 0 \quad \left\{ \begin{array}{l} 1 + 1 \geq \gamma \lambda_i \\ 2 \geq \gamma \lambda_i \end{array} \right. \quad \left. \begin{array}{l} 2 \geq \gamma \lambda_i \\ \lambda_i \leq \frac{2}{\gamma} \quad \forall i \end{array} \right\}$ use ②

• the evolution for B is stable if
 $0 \leq \lambda_i \leq \frac{2}{\gamma} \quad \forall$ eigenvalue λ_i of $A^T A$

$$|U^T V| \leq \|U\| \|V\|$$

3.c. $\varphi(x) = x - \gamma \nabla F(x)$ $\gamma > 0$ constant step size

Show $\forall x, x' \in \mathbb{R}^d \quad \|\varphi(x) - \varphi(x')\|_2 \leq \beta \|x - x'\|_2$

$$\beta = \max\{|1 - \gamma \lambda_{\max}(A^T A)|, |1 - \gamma \lambda_{\min}(A^T A)|\}$$

$$\|x - \gamma \nabla F(x) - x' + \gamma \nabla F(x')\|_2 = \|x - x' + \gamma (\nabla F(x') - \nabla F(x))\|_2$$

$$= \|x - x' + \gamma (A^T A x' - A^T b - A^T A x + A^T b)\|_2 \quad \text{apply rule for gradients}$$

$$= \|(I - \gamma A^T A)(x - x')\|_2$$

$$= \|(I - \gamma A^T A)(x - x')\|_2 \leq \|I - \gamma A^T A\| \|x - x'\|_2$$

$$\|I - \gamma A^T A\| = \sqrt{\lambda_{\max}(I - \gamma A^T A)} = \|\gamma \lambda_{\max}(A^T A)\| \quad \text{• always pos}$$

* From matrix norm with defn

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A) \quad \text{• } \sigma \text{ must always be positive}$$

From 3.b, the eigen values of $I - \gamma A^T A$ are bounded by $0 \leq \lambda \leq \frac{\gamma}{2}$

$$\|I - \gamma A^T A\| = \max(|1 - \gamma \lambda_{\max}(A^T A)|) = \max\{|1 - \gamma \lambda_{\max}(A^T A)|, |1 - \gamma \lambda_{\min}(A^T A)|\}$$

for all eigs

$$= \beta$$

$$\leq \beta \|x - x'\|_2 \text{ as required}$$

$$3.c. \quad x^* = \arg \min_{x \in \mathbb{R}^n} f(x) , \quad \varphi(x) = x - \gamma \nabla f(x)$$

$$\text{Show: } \|x_{k+1} - x^*\|_2 = \|\varphi(x_k) - \varphi(x^*)\|_2$$

$$\text{From 3.c. we know } \|\varphi(x_k) - \varphi(x^*)\| \leq \beta \|x_k - x^*\|$$

$$\left. \begin{array}{l} \varphi(x^*) = x^* - \gamma \nabla f(x^*) \\ \varphi(x) = x_{k+1} = x_k - \gamma \nabla f(x_k) \end{array} \right\} \therefore \|x_{k+1} - x^*\|_2 = \|\varphi(x_k) - \varphi(x^*)\|_2$$

$$\text{conclude } \|x_{k+1} - x^*\|_2 \leq \beta^{k+1} \|x_0 - x^*\|_2$$

but (from 3.c)

$$\|x_{k+1} - x^*\|_2 \leq \beta \|x_k - x^*\|_2$$

$$\|x_1 - x^*\|_2 \leq \beta \|x_0 - x^*\|_2 \quad \text{- from 3.c.}$$

$$\begin{aligned} \|x_2 - x^*\|_2 &\leq \beta \|x_1 - x^*\|_2 \leq \beta (\beta \|x_0 - x^*\|_2) \\ &\leq \beta^2 \|x_0 - x^*\|_2 \end{aligned}$$

$$\|x_{k+1} - x^*\|_2 \leq \beta^{k+1} \|x_0 - x^*\|_2 \quad \text{as required}$$

$$3.e \quad \textcircled{1} \quad f(x) - f(x^*) = \frac{1}{2} \|Ax - x^*\|^2 \quad \textcircled{2} \quad f(x) = \frac{1}{2} \|Ax - b\|^2$$

$Ax = x^* \quad (Ax - Ax^*)$

$$\frac{1}{2} \left(\|Ax - b\|^2 - \|Ax^* - b\|^2 \right) = \frac{1}{2} (x - x^*)^T A^T A (x - x^*)$$

$$\frac{1}{2} \left[x^T A^T A x - x^T A^T b - b^T A x + b^T b = x^T A^T A x^* + x^T A^T b + b^T A x^* - b^T b \right] \\ - 2b^T A x + 2b^T A x^*$$

x^* is opt. soln for $\frac{1}{2} \|Ax - b\|^2 \quad x^* = (A^T A)^{-1} A^T b + A^T A \text{pos}$

$$\frac{1}{2} \left(x^T A^T A x - 2b^T A x + 2b^T A (A^T A)^{-1} A^T b = x^*^T A^T A x^* \right)$$

$\text{RHS} = b^T b$

$$(A^T A)^{-1} A^T A (A^T A)^{-1} A^T b = -b^T b$$

$$2b^T A (A^T A)^{-1} A^T b = 2b^T A A^{-1} A^T A^T b = 2b^T b$$

$$= \frac{1}{2} (x^T A^T A x - 2b^T A x + 2b^T b) = \frac{1}{2} (x^T A^T A x - 2b^T A x + b^T b) \quad \textcircled{1}$$

$$\textcircled{2} \quad \frac{1}{2} (x^T A^T A x - x^* T A^T A x^* - x^* T A^T A x + x^* T A^T A x)$$

$$= \frac{1}{2} (x^T A^T A x - x^* T A^T A (A^T A)^{-1} A^T b - [(A^T A)^{-1} A^T b]^T A^T A x + b^T b)$$

$- x^* T A^T b \quad b^T A (A^T A)^{-1} A^T A x$

$$= \frac{1}{2} (x^T A^T A x - x^* T A^T b - b^T A x + b^T b) = \frac{1}{2} (x^T A^T A x - 2b^T A x + b^T b) \quad \textcircled{2}$$

LHS = RHS

\textcircled{1}

$= (b^T A (A^T A)^{-1}) A^T A A^T A^{-1} A^T b = -b^T b$
--

\textcircled{2}

$$3.F. \quad F(x_k) - F(x^*) \leq \frac{\kappa}{2} \|x_k - x^*\|_2^2 \quad \kappa = 2\max(A^T A)$$

Show the above statement:

$$\text{From 3.e. } F(x_k) - F(x^*) = \frac{1}{2} \|A(x_k - x^*)\|_2^2 :$$

$$\begin{aligned} \|A\|_2^2 &= (\sqrt{2\max(A^T A)})^2 = 2\max(A^T A) = \kappa \\ &\leq \frac{\kappa}{2} \|x_k - x^*\|_2^2 \end{aligned}$$

From 3.d. we conclude the following:

$$\frac{\kappa}{2} \|x_k - x^*\|_2^2 \leq \frac{\kappa}{2} \beta^{2k} \|x_0 - x^*\|_2^2$$

$$\therefore F(x_k) - F(x^*) \leq \frac{\kappa}{2} \beta^{2k} \|x_0 - x^*\|_2^2 \text{ as required}$$

$$3.g. \quad \beta \text{ is function of } \gamma \quad \kappa = \frac{2\max(A^T A)}{2\min(A^T A)}$$

$$\text{Effectively pick } \gamma \text{ s.t. } \beta = \frac{\kappa-1}{\kappa+1} = \frac{2\max(A^T A)}{2\min(A^T A)} - 1$$

$$\therefore \frac{2\max(A^T A) - 2\min(A^T A)}{2\min(A^T A)}, \quad \frac{2\max(A^T A)}{2\min(A^T A)} + 1$$

$$= \frac{2\max(A^T A) - 2\min(A^T A)}{2\max(A^T A) + 2\min(A^T A)} = \frac{1 - 2\min(A^T A)}{1 + 2\min(A^T A)} = \max\left\{ \left| 1 - 2\min(A^T A) \right|, \left| 1 + 2\min(A^T A) \right| \right\}$$

(cont)

$$\text{3.g. } 1 - \gamma \lambda_{\max}(A^T A) = \frac{\lambda_{\max}(A^T A) - \lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A) + \lambda_{\min}(A^T A)}$$

$$\lambda_{\max}(A^T A) = M$$

$$\lambda_{\min}(A^T A) = N$$

$$1 - \gamma M = \frac{M - N}{M + N}, \quad 1 - \frac{(M - N)}{M + N} = \gamma M$$

$$\frac{\gamma N}{M + N}$$

$$\frac{M + N - M + N}{M + N} = \gamma M, \quad \left(\gamma = \frac{2N}{(M + N)M} \right) \quad \text{for } 1 - \lambda_{\max}(A^T A)$$

$$\gamma = \frac{2 \lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A) \lambda_{\min}(A^T A) + \lambda_{\min}(A^T A) \cdot \lambda_{\max}(A^T A)}$$

same

v

$$\gamma = \frac{2N}{(M + N)N}$$

for $1 - \lambda_{\min}(A^T A)$

$$\frac{2N}{MN + NN}$$

$$\gamma = \frac{2 \lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A) \lambda_{\min}(A^T A) + \lambda_{\min}(A^T A) \lambda_{\min}(A^T A)}$$

• γ can be either of the above. $\beta = \max \{ |1 - \gamma \lambda_{\max}(A^T A)|, |1 - \gamma \lambda_{\min}(A^T A)| \}$

• want γ & β as small as possible

Muscles, Nerves Distance: Dij
(alpha) (beta) (chi)

4.a. $M=2$, $\Lambda=1$

$$D_{ij} = \|(\mathbf{a}_i, \mathbf{b}_i) - (\mathbf{x}_j, \mathbf{y}_j)\| + \epsilon_{ij} \sim N(0, 1)$$

$$L(x, y_i) = - \sum_{j=1}^3 \left(\left((a_j - x_j)^2 + (b_j - y_j)^2 \right)^{1/2} - d_j \right)^2$$

$$\nabla_{x,y} \left\{ \delta L(x,y) / s_x = -2 \sum_{i=1}^n \left((a_i - x_i)^2 + (b_i - y_i)^2 \right)^{1/2} - c_i \right)$$

$$+ \frac{1}{2} \left((a_i - x_i)^2 + (b_i - y_i)^2 \right)^{-1/2} \cdot 2(a_i - x_i) \cdot (-1)$$

$$\delta L(x,y) / s_y = -2 \sum_{i=1}^n \left((a_i - x_i)^2 + (b_i - y_i)^2 \right)^{1/2} - c_i \cdot \frac{1}{2} \left((a_i - x_i)^2 + (b_i - y_i)^2 \right)^{-1/2}$$

$$\cdot 2(b_i - y_i) \cdot (-1)$$

$$d_{x,y} = \sqrt{2 \sum_{i=1}^n ((a_i - x_i)^2 + (b_i - y_i)^2)^{1/2} - d_i} / \left(\sum_{i=1}^n ((a_i - x_i)^2 + (b_i - y_i)^2)^{-1/2} \right)$$

$$S. f(x) = \sin(x)$$

$$a_{i+1} = f_i(a_i) = \sigma(z_i) = \sigma(w_i a_i + b_i)$$

$$a_{i+1} =$$

weights

$W \in \mathbb{R}^{N_i \times N_{i+1}}$

N_{i+1} : size of layer $i+1$

$$\frac{\partial f_i}{\partial a_i} =$$

b_i : vector biases

$$M_i = N_{i-1}$$

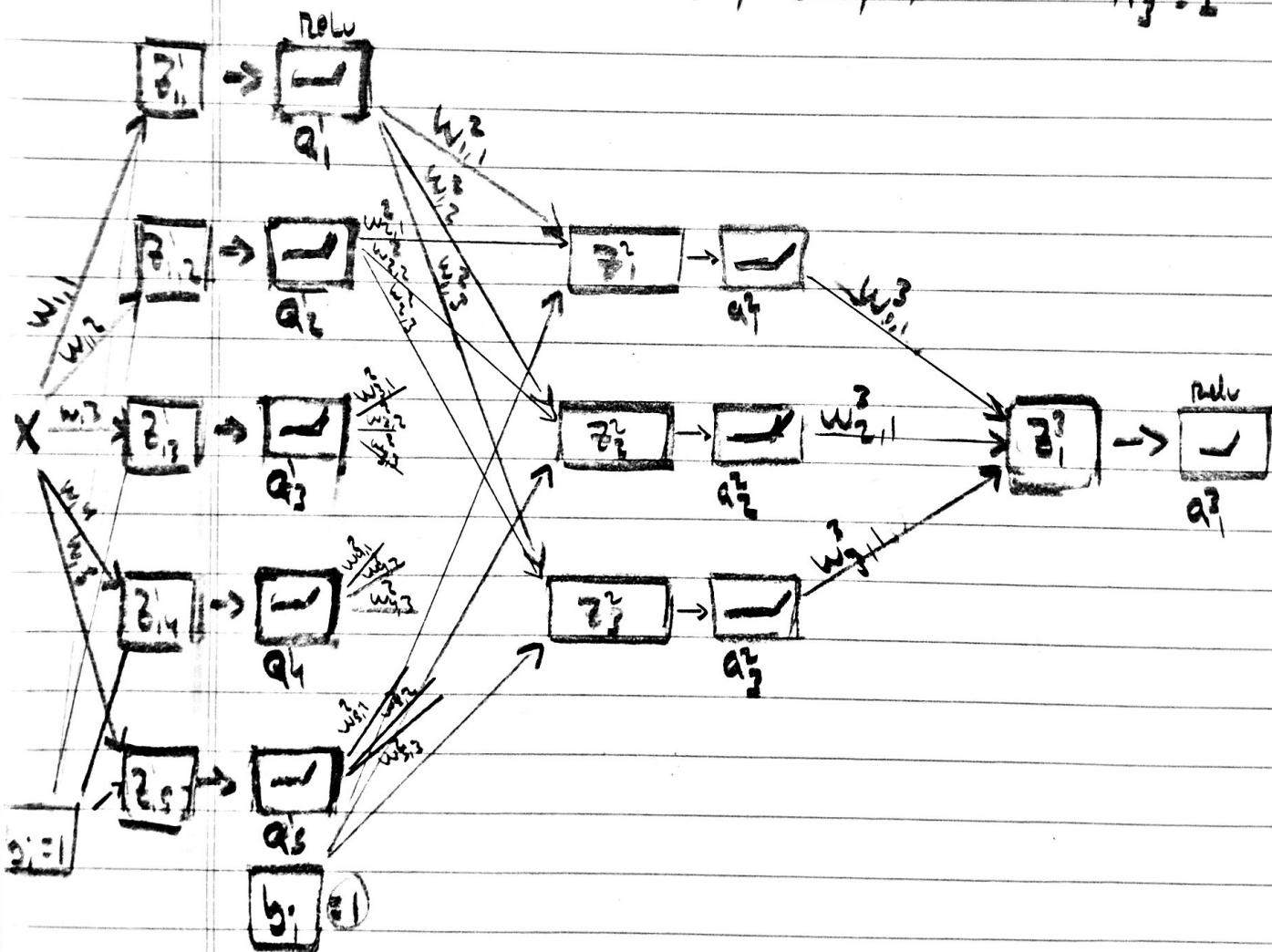
$$a_i \in \mathbb{R}^{M_i \times 1}$$

$$w_a \rightarrow \mathbb{R}^{N_i \times 2}$$

$$N_i = 3$$

$$z_i \\ (w_i a_i + b_i) \mathbb{R}^{N_i \times 2}$$

a. example network w/ 3 computational layers $N_1 = 5$ $b_1, b_2, b_3 \neq 1$
 $N_2 = 3$
 $N_3 = 2$



key $[w_{i,j}]$ w is the layer i
 i - from i to j about]

$$S.b \quad MSE(\hat{y}) = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad f(x) = \sin(x) = y_i$$

$$\hat{y} = h(x)$$

$$\begin{aligned} \frac{\partial (MSE(\hat{y}))}{\partial \hat{y}} &= \frac{1}{2} \sum_{i=1}^n 2(y_i - \hat{y}_i)^1 \cdot (-1) \cdot \frac{\partial \hat{y}_i}{\partial \hat{y}} \\ &= - \sum_{i=1}^n (y_i - \hat{y}_i) \cdot 1, \quad \boxed{= (\hat{y}_i - y_i)} \end{aligned}$$

$$S.c. \quad \sigma_{\text{tanh}}(z) = ? \quad \cdot \frac{\partial \sigma_z}{\partial z} = ?$$

$$\cdot \sigma_{\text{tanh}}(z) = (e^z - e^{-z})(e^z + e^{-z})^{-1} =$$

$$\begin{aligned} \cdot \frac{\partial \sigma_{\text{tanh}}}{\partial z} &= (e^z + e^{-z})(e^z + e^{-z})^{-1} + (e^z - e^{-z})(-1)(e^z + e^{-z})^{-2} \\ &= 1 - (e^z - e^{-z})(e^z + e^{-z})^{-2} / (e^z - e^{-z}) \end{aligned}$$

$$= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$\cdot \sigma_{\text{ReLU}}(z) = \begin{cases} 0 & z < 0 \\ 1 & \text{o.w.} \end{cases} \quad \cdot \frac{\partial \sigma_{\text{ReLU}}}{\partial z} = \begin{cases} 0 & z < 0 \\ 1 & \text{o.w.} \end{cases}$$

in code key want
element wise sum

S.d. Given: $\frac{\delta \text{MSE}}{\delta a_{i+1}}$ a_{i+1} is input to layer $i+1$

$$\frac{\delta \text{MSE}}{\delta a_i} = \left(\text{in terms of } \frac{\delta \text{MSE}}{\delta a_{i+1}} \right) \text{ or } \left| \frac{\delta \text{MSE}}{\delta a_i} \right|$$

$$a_{i+1} = \sigma(w_i a_i + b_i), \hat{a}_{i+1} = \underbrace{\sigma(w_i a_i + b_i)}_{\text{output}}$$

$$\text{MSE}(\hat{a}_i) = \frac{1}{2} \sum (a_i - \hat{a}_i)^2$$

Bias prop so $a_{i+1} \rightarrow a_i'$

$$\text{MSE}(\hat{a}_{i+1}) = \frac{1}{2} (a_{i+1} - \hat{a}_{i+1})^2 = \frac{1}{2} (a_{i+1} - \sigma(w_i a_i + b_i))^2$$

$$\text{MSE}(\hat{a}_i) = \frac{1}{2} (a_i - \hat{a}_i)^2 = \frac{1}{2} (a_i - \sigma(w_i a_{i-1} + b_{i-1}))^2$$

$$\frac{\delta \text{MSE}}{\delta a_{i+1}} = \frac{\delta \text{MSE}}{\delta a_i}$$

$$\frac{\delta \text{MSE}}{\delta a_i} = \frac{1}{2} \frac{\delta (a_i - \hat{a}_i)^2}{\delta a_{i+1}} \cdot \frac{\delta a_{i+1}}{\delta a_i} =$$

$$\frac{\delta a_{i+1}}{\delta a_i} = \frac{\delta (\sigma(z_i))}{\delta z_i} \cdot \frac{\delta z_i}{\delta a_i} \quad z_i = w_i a_i + b_i$$

$$\frac{\delta \text{MSE}}{\delta a_i} = \frac{\delta \text{MSE}}{\delta a_{i+1}} \cdot \frac{\delta (\sigma(z_i))}{\delta z_i} \cdot \frac{\delta z_i}{\delta a_i}$$

$$\frac{\delta a_{i+1}}{\delta a_i}$$

this is w_i

Question 2 B

In [39]:

```
""" Tools for calculating Gradient Descent for ||Ax-b||. """
import matplotlib.pyplot as plt
import numpy as np


def compute_gradient(A, b, x):
    """Computes the gradient of ||Ax-b|| with respect to x."""
    return np.dot(A.T, (np.dot(A, x) - b)) / np.linalg.norm(np.dot(A, x) - b)

def compute_update(A, b, x, step_count, step_size):
    """Computes the new point after the update at x."""
    return x - step_size(step_count) * compute_gradient(A, b, x)

def compute_updates(A, b, p, total_step_count, step_size):
    """Computes several updates towards the minimum of ||Ax-b|| from p.

    Params:
        b: in the equation ||Ax-b||
        p: initialization point
        total_step_count: number of iterations to calculate
        step_size: function for determining the step size at step i
    """
    positions = [np.array(p)]
    n = 0
    for k in range(total_step_count):
        n += 1
        positions.append(compute_update(A, b, positions[-1], k, step_size))
    print(n)
    return np.array(positions)

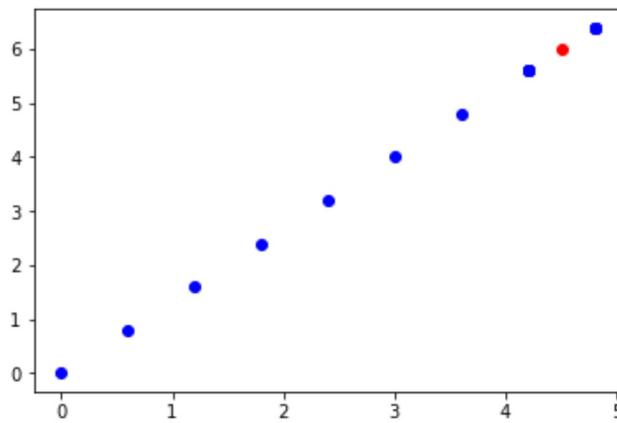
#####
# TODO(student): Input Variables
A = np.array([[1, 0], [0, 1]]) # do not change this until the last part
b = np.array([4.5, 6]) # b in the equation ||Ax-b||
initial_position = np.array([0, 0]) # position at iteration 0
total_step_count = 15 # number of GD steps to take
step_size = lambda i: 1 # step size at iteration i
#####

# computes desired number of steps of gradient descent
positions = compute_updates(A, b, initial_position, total_step_count, step_size)

# print out the values of the x_i
print(positions)
print(np.dot(np.linalg.inv(A), b))

# plot the values of the x_i
plt.scatter(positions[:, 0], positions[:, 1], c='blue')
plt.scatter(np.dot(np.linalg.inv(A), b)[0],
            np.dot(np.linalg.inv(A), b)[1], c='red')
plt.plot()
plt.show()
```

```
15  
[[ 0.   0. ]  
 [ 0.6  0.8]  
 [ 1.2  1.6]  
 [ 1.8  2.4]  
 [ 2.4  3.2]  
 [ 3.   4. ]  
 [ 3.6  4.8]  
 [ 4.2  5.6]  
 [ 4.8  6.4]  
 [ 4.2  5.6]  
 [ 4.8  6.4]  
 [ 4.2  5.6]  
 [ 4.8  6.4]  
 [ 4.2  5.6]  
 [ 4.5  6. ]]
```



Although the function is convex, the gradient descent method will not find the optimal solution.

It will never get within 0.01 of the optimal solution. The gradient descent method clearly begins oscillating between [4.8, 6.4] and [4.2, 5.6]. The error in these steps are [0.3, 0.4], which does not satisfy epsilon. Does not converge to the optimal solution.

The step size is constant and is too large.

For General b

In [54]:

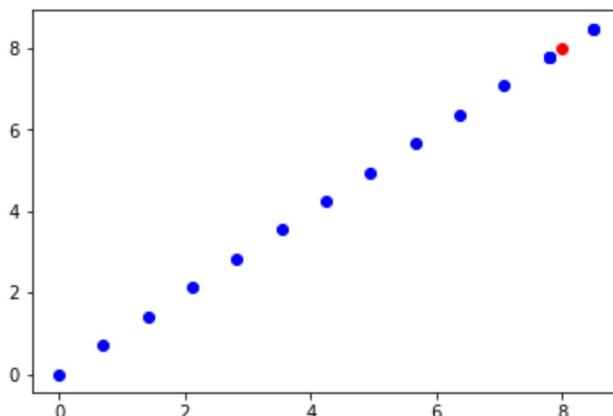
```
#####
# TODO(student): Input Variables
A = np.array([[1, 0], [0, 1]]) # do not change this until the last part
b = np.array([8, 8]) # b in the equation ||Ax-b||
initial_position = np.array([0, 0]) # position at iteration 0
total_step_count = 15 # number of GD steps to take
step_size = lambda i: 1 # step size at iteration i
#####
#
# computes desired number of steps of gradient descent
positions = compute_updates(A, b, initial_position, total_step_count, step_size)
)

# print out the values of the x_i
print(positions)
print(np.dot(np.linalg.inv(A), b))

# plot the values of the x_i
plt.scatter(positions[:, 0], positions[:, 1], c='blue')
plt.scatter(np.dot(np.linalg.inv(A), b)[0],
            np.dot(np.linalg.inv(A), b)[1], c='red')
plt.plot()
plt.show()
```

15

```
[[ 0.          0.          ]
 [ 0.70710678  0.70710678]
 [ 1.41421356  1.41421356]
 [ 2.12132034  2.12132034]
 [ 2.82842712  2.82842712]
 [ 3.53553391  3.53553391]
 [ 4.24264069  4.24264069]
 [ 4.94974747  4.94974747]
 [ 5.65685425  5.65685425]
 [ 6.36396103  6.36396103]
 [ 7.07106781  7.07106781]
 [ 7.77817459  7.77817459]
 [ 8.48528137  8.48528137]
 [ 7.77817459  7.77817459]
 [ 8.48528137  8.48528137]
 [ 7.77817459  7.77817459]
 [ 8.        8.]]
```



The same can be said for general b. Gradient descent will not always reach the optimal solution.

The step size is constant and is too large. Therefore it will not always converge.

Question 2 C

In [60]:

```
#####
# TODO(student): Input Variables
A = np.array([[1, 0], [0, 1]]) # do not change this until the last part
b = np.array([4.5, 6]) # b in the equation ||Ax-b||
initial_position = np.array([0, 0]) # position at iteration 0
total_step_count = 5000 # number of GD steps to take
step_size = lambda i: (5/6)**i # step size at iteration i
#####
# computes desired number of steps of gradient descent
positions = compute_updates(A, b, initial_position, total_step_count, step_size)

# print out the values of the x_i
print(positions)
print(np.dot(np.linalg.inv(A), b))

# plot the values of the x_i
plt.scatter(positions[:, 0], positions[:, 1], c='blue')
plt.scatter(np.dot(np.linalg.inv(A), b)[0],
            np.dot(np.linalg.inv(A), b)[1], c='red')
plt.plot()
plt.show()

# computes desired number of steps of gradient descent
b = [3, 3]
positions = compute_updates(A, b, initial_position, total_step_count, step_size)

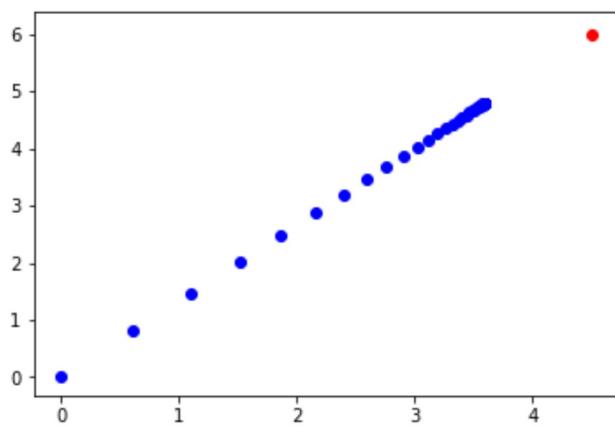
# print out the values of the x_i
print(positions)
print(np.dot(np.linalg.inv(A), b))

# plot the values of the x_i
plt.scatter(positions[:, 0], positions[:, 1], c='blue')
plt.scatter(np.dot(np.linalg.inv(A), b)[0],
            np.dot(np.linalg.inv(A), b)[1], c='red')
plt.plot()
plt.show()
```

```

5000
[[ 0.          0.          ]
 [ 0.6         0.8         ]
 [ 1.1         1.46666667]
 ...,
 [ 3.6         4.8         ]
 [ 3.6         4.8         ]
 [ 3.6         4.8         ]]
[ 4.5   6. ]

```



```

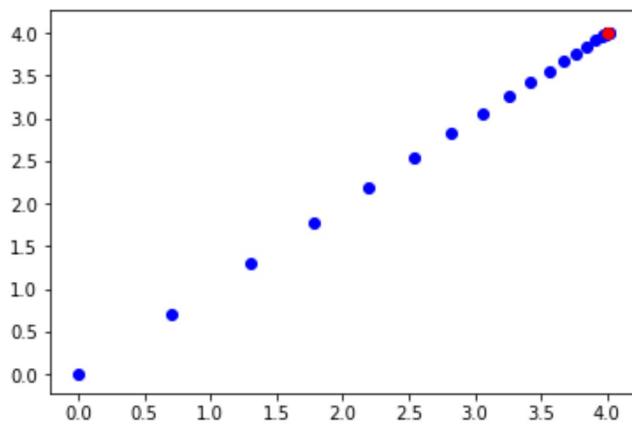
/opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:10: RuntimeWarning:
invalid value encountered in true_divide
# Remove the CWD from sys.path while we load stuff.

```

```

5000
[[ 0.          0.          ]
 [ 0.70710678  0.70710678]
 [ 1.29636243  1.29636243]
 ...,
 [      nan      nan]
 [      nan      nan]
 [      nan      nan]]
[ 4.   4.]

```



No the gradient descent will not find the optimal solution. This is because our step size gets exponentially smaller as we iterate. It does not converge. The plot does a good job illustrating this point.

For General b

Conceptually, if our optimal solution x^* is close enough to our initial x than our gradient descent method will converge to the optimal solution.

Question 2 D

```
In [57]: ##### Input Variables
# TODO(student): Input Variables
A = np.array([[1, 0], [0, 1]]) # do not change this until the last part
b = np.array([4.5, 6]) # b in the equation ||Ax-b||
initial_position = np.array([0, 0]) # position at iteration 0
total_step_count = 5000 # number of GD steps to take
step_size = lambda i: 1/(i + 1) # step size at iteration i
##### Compute desired number of steps of gradient descent
positions = compute_updates(A, b, initial_position, total_step_count, step_size)

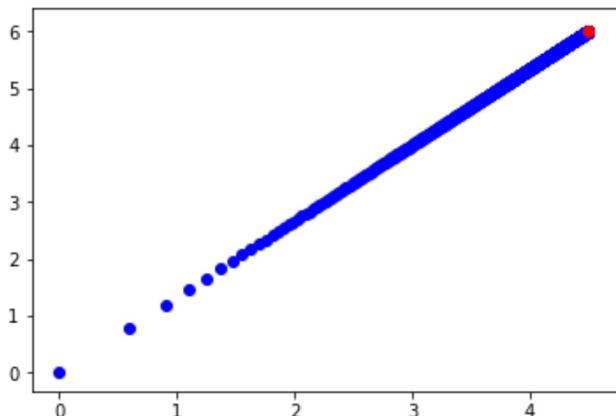
# print out the values of the x_i
print(positions)
print(np.dot(np.linalg.inv(A), b))

# plot the values of the x_i
plt.scatter(positions[:, 0], positions[:, 1], c='blue')
plt.scatter(np.dot(np.linalg.inv(A), b)[0],
            np.dot(np.linalg.inv(A), b)[1], c='red')
plt.plot()
plt.show()
```

```

5000
[[ 0.           0.          ]
 [ 0.6          0.8         ]
 [ 0.9          1.2         ]
...
[ 4.50012003  6.00016004]
[ 4.50000001  6.00000001]
[ 4.49988001  5.99984001]]
[ 4.5  6. ]

```



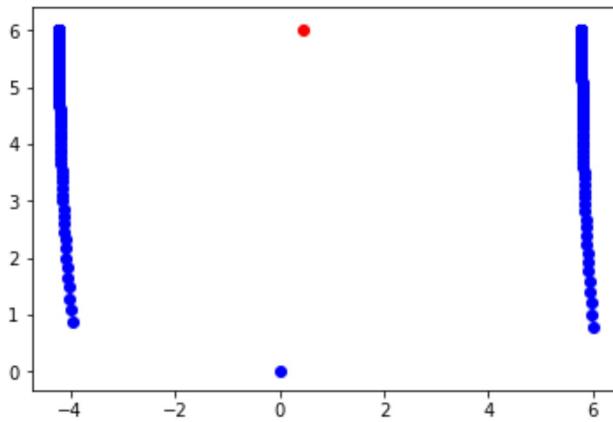
Our gradient descent method converges to the optimal solution. The expression for the number of iterations required for convergence is in the hand written part of this HW submission.

This step size makes our gradient descent method work for general b .

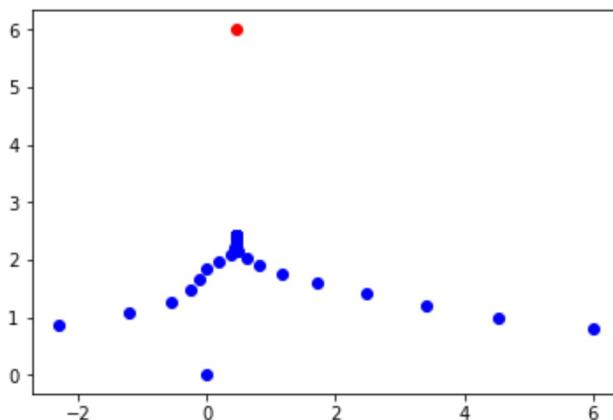
Question 2 E

```
In [71]: #####  
# TODO(student): Input Variables  
A = np.array([[10, 0], [0, 1]]) # do not change this until the last part  
b = np.array([4.5, 6]) # b in the equation ||Ax-b||  
initial_position = np.array([0, 0]) # position at iteration 0  
total_step_count = 5000 # number of GD steps to take  
steps = [lambda i : 1, lambda i: (5/6)**i, lambda i: 1/(i + 1)] # step size at iteration i  
#####  
  
# plot the values of the x_i  
for i in range(3):  
    print(i)  
    step_size = steps[i]  
    # computes desired number of steps of gradient descent  
    positions = compute_updates(A, b, initial_position, total_step_count, step_size)  
    # print out the values of the x_i  
    print(positions)  
    print(np.dot(np.linalg.inv(A), b))  
    plt.scatter(positions[:, 0], positions[:, 1], c='blue')  
    plt.scatter(np.dot(np.linalg.inv(A), b)[0],  
                np.dot(np.linalg.inv(A), b)[1], c='red')  
    print(A, "Step Size Method: " + str(i))  
    plt.plot()  
    plt.show()  
  
#####  
# TODO(student): Input Variables  
A = np.array([[15, 8], [6, 5]]) # do not change this until the last part  
#####  
  
# plot the values of the x_i  
for i in range(3):  
    print(i)  
    step_size = steps[i]  
    # computes desired number of steps of gradient descent  
    positions = compute_updates(A, b, initial_position, total_step_count, step_size)  
    # print out the values of the x_i  
    print(positions)  
    print(np.dot(np.linalg.inv(A), b))  
    plt.scatter(positions[:, 0], positions[:, 1], c='blue')  
    plt.scatter(np.dot(np.linalg.inv(A), b)[0],  
                np.dot(np.linalg.inv(A), b)[1], c='red')  
    print(A, "Step Size Method: " + str(i))  
    plt.plot()  
    plt.show()
```

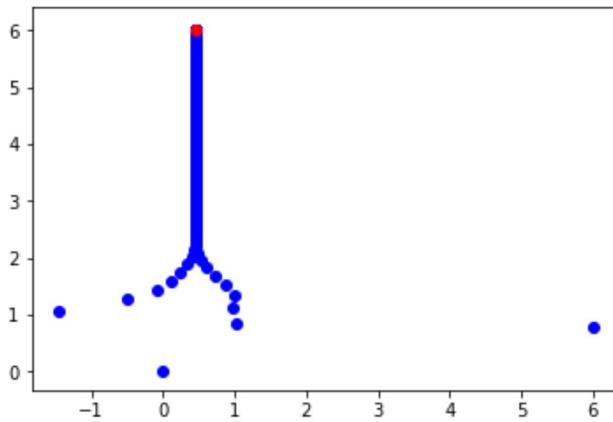
```
0
5000
[[ 0.          0.          ]
 [ 6.          0.8         ]
[-3.95639434  0.89328514]
....,
[-4.21969358  6.          ]
[ 5.78030642  6.          ]
[-4.21969358  6.          ]]
[ 0.45  6.  ]
[[10  0]
 [ 0  1]] Step Size Method: 0
```



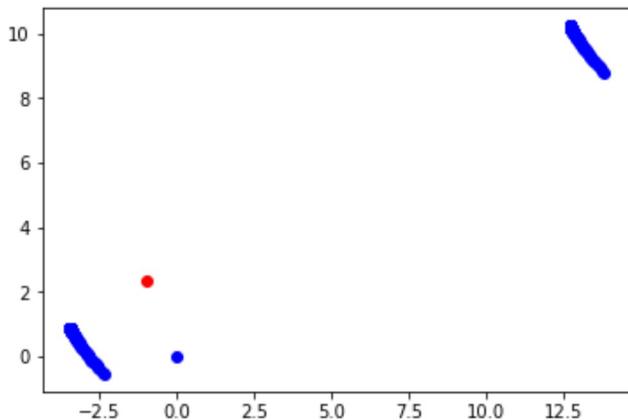
```
1
5000
[[ 0.          0.          ]
 [ 6.          0.8         ]
[-2.29699529  0.87773761]
....,
[ 0.44999997  2.43095698]
[ 0.44999997  2.43095698]
[ 0.44999997  2.43095698]]
[ 0.45  6.  ]
[[10  0]
 [ 0  1]] Step Size Method: 1
```



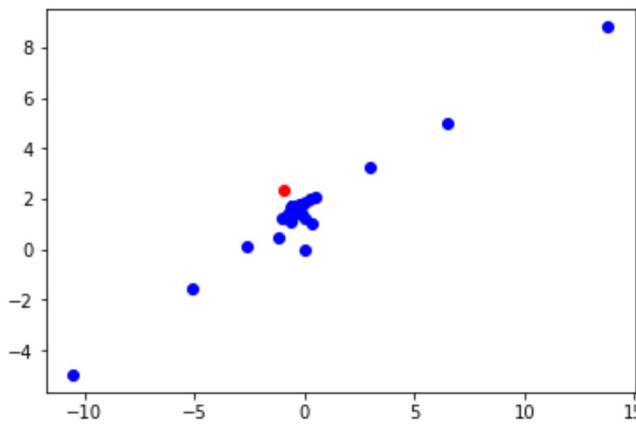
```
2
5000
[[ 0.          0.          ]
 [ 6.          0.8         ]
 [ 1.02180283  0.84664257]
 ....,
 [ 0.45        6.00020008]
 [ 0.45        6.00000004]
 [ 0.45        5.99980004]]
[ 0.45  6.  ]
[[10  0]
 [ 0  1]] Step Size Method: 2
```



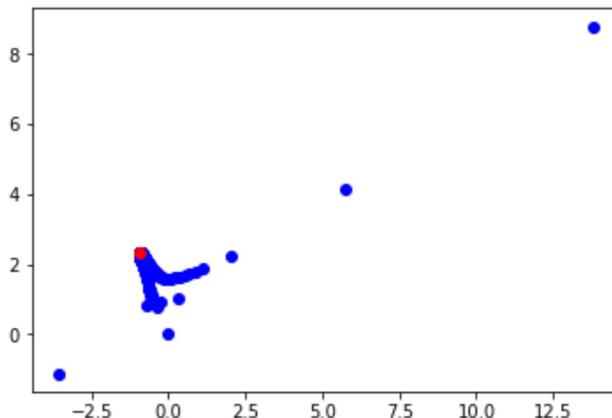
```
0
5000
[[ 0.          0.          ]
 [ 13.8        8.8         ]
 [ -2.34489789 -0.53919835]
 ....,
 [ -3.45081002  0.88123538]
 [ 12.68837859  10.23170003]
 [ -3.45081002  0.88123538]]
[-0.94444444  2.33333333]
[[15  8]
 [ 6  5]] Step Size Method: 0
```



```
1
5000
[[ 0.          0.          ]
 [ 13.8        8.8         ]
 [  0.34591842  1.01733471]
 ...,
 [ -0.58650742  1.715523  ]
 [ -0.58650742  1.715523  ]
 [ -0.58650742  1.715523  ]]
[-0.94444444  2.33333333]
[[15  8]
 [ 6  5]] Step Size Method: 1
```



```
2
5000
[[ 0.          0.          ]
 [ 13.8        8.8         ]
 [  5.72755105  4.13040083]
 ...,
 [ -0.94283004  2.33426866]
 [ -0.94605852  2.33239819]
 [ -0.94283069  2.33426829]]
[-0.94444444  2.33333333]
[[15  8]
 [ 6  5]] Step Size Method: 2
```



The step size function that we utilized for Q2.d should converge to the optimal solution regardless of the matrix A.

Changing A only changes what we are minimizing.

Question 5

```
In [98]: import numpy as np
import matplotlib
import matplotlib.pyplot as plt

# Gradient descent optimization
# The learning rate is specified by eta
class GDOptimizer(object):
    def __init__(self, eta):
        self.eta = eta

    def initialize(self, layers):
        pass

    # This function performs one gradient descent step
    # layers is a list of dense layers in the network
    # g is a list of gradients going into each layer before the nonlinear activation
    # a is a list of activations of each node in the previous layer going
    #
    def update(self, layers, g, a):
        m = a[0].shape[1]
        for layer, curGrad, curA in zip(layers, g, a):
            # TODO: PART F #####
            # Compute the gradients for layer.W and layer.b using the gradient for
            # the output of the
            # layer curA and the gradient of the output curGrad
            # Use the gradients to update the weight and the bias for the layer
            #
            # Normalize the learning rate by m (defined above), the number of training examples input
            # (in parallel) to the network.
            #
            # It may help to think about how you would calculate the update if we input just one
            # training example at a time; then compute a mean over these individual update values.
            # #####
            Wgrad_layer = curGrad * curA
            bgrad_layer = curGrad

            layer.W = layer.W - self.eta/m * np.mean(Wgrad_layer)
            layer.b = layer.b - self.eta/m * np.mean(bgrad_layer)

    # Cost function used to compute prediction errors
    class QuadraticCost(object):

        # Compute the squared error between the prediction yp and the observation y
        # This method should compute the cost per element such that the output is the
        # same shape as y and yp
        @staticmethod
        def fx(y,yp):
            # TODO: PART B #####
            # Implement me
            return 0.5*np.square(y - yp)
            # #####
            # Derivative of the cost function with respect to yp
        @staticmethod
        def dx(v, v0):
```

Debugging gradients..

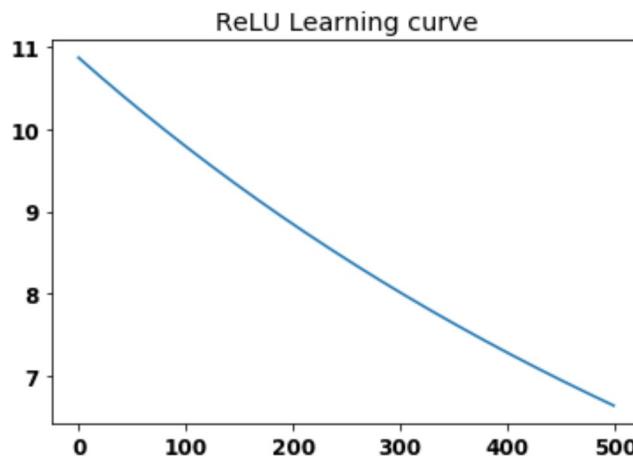
squared difference of layer 0: 3.05989714594e-10

squared difference of layer 1: 1.88645317382e-10

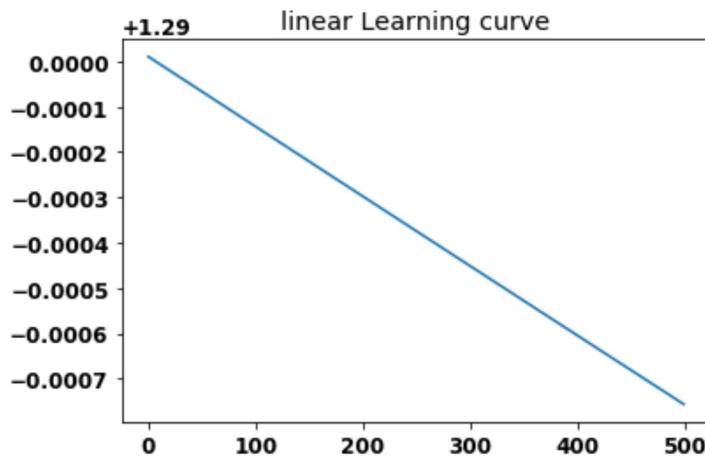
squared difference of layer 2: 9.84497101272e-11

Standard fully connected network

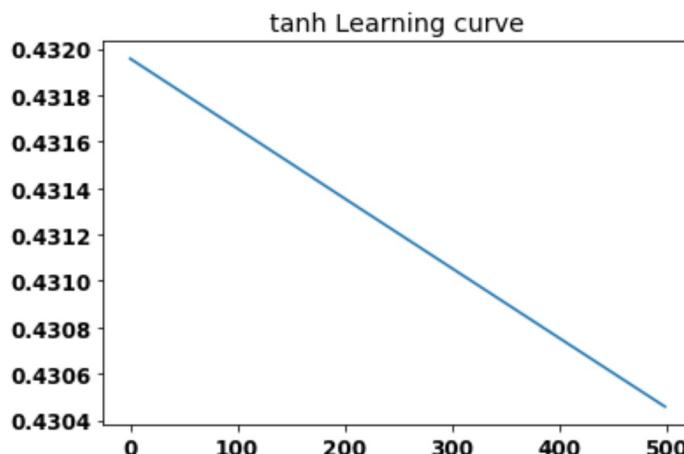
ReLU MSE: 6.64656178094

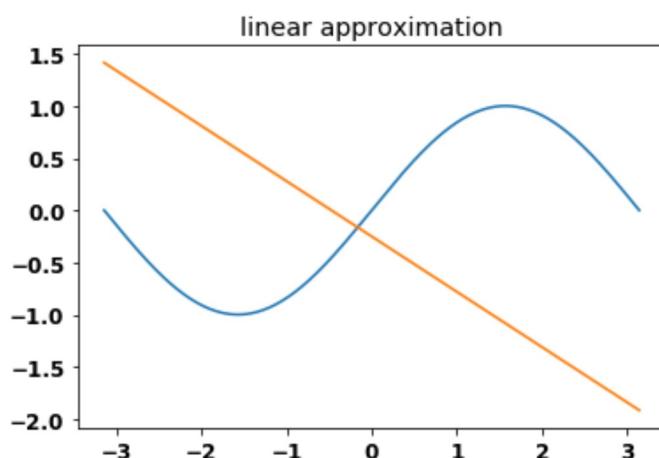
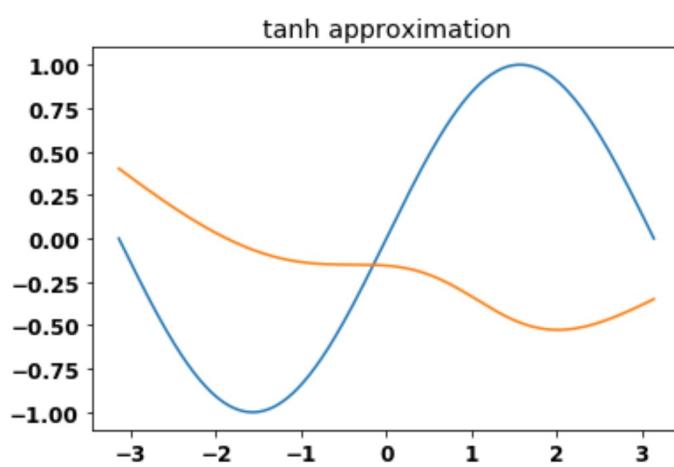
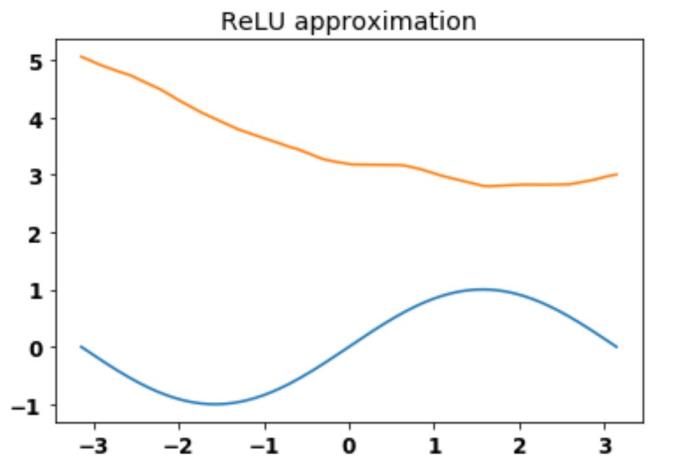


linear MSE: 1.28924330257



tanh MSE: 0.430458640433





Training with various sized network

ReLU MSE Error

Layers	5 nodes	10 nodes	25 nodes	50 nodes
1	1.90620	1.64771	0.32309	0.13785
2	0.28166	0.73836	2.74120	1.16850
3	5.60639	6.09030	5.19953	2.26372

tanh MSE Error

Layers	5 nodes	10 nodes	25 nodes	50 nodes
1	0.68136	2.37873	4.06720	0.25491
2	1.00403	0.10372	0.48874	3.52292
3	0.27115	1.76387	0.15810	1.71457

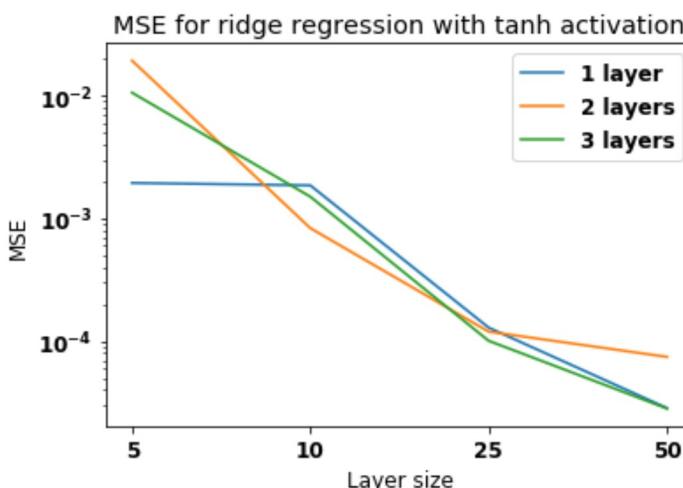
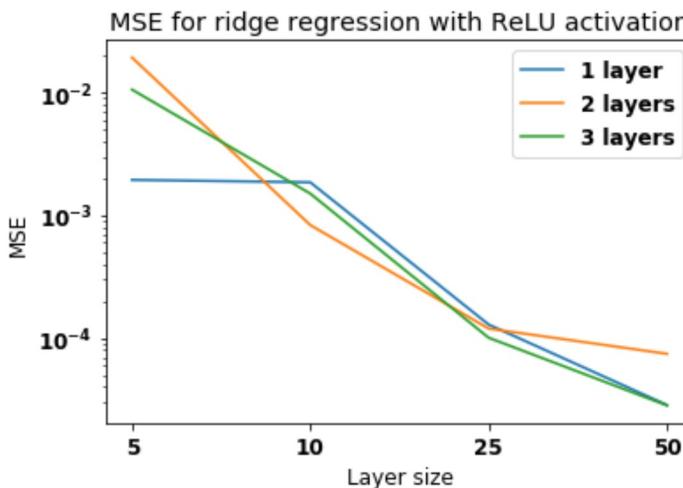
Running ridge regression on last layer

ReLU MSE Error

Layers	5 nodes	10 nodes	25 nodes	50 nodes
1	0.00696	0.01074	0.00067	0.00005
2	0.01916	0.00284	0.00018	0.00007
3	0.03777	0.00106	0.00009	0.00009

tanh MSE Error

Layers	5 nodes	10 nodes	25 nodes	50 nodes
1	0.00195	0.00186	0.00013	0.00003
2	0.01931	0.00083	0.00012	0.00007
3	0.01058	0.00150	0.00010	0.00003



Question 5 Comment

Layer size is proportional to MSE. For both nonlinearities the MSE decreases with layer size. However there is not an appreciable difference as we increase the number of layers.

```
In [7]: """ Tools for calculating Gradient Descent for ||Ax-b||. """
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
```

4. Sensors, Objects, and Localization

```
In [8]: import numpy as np
import scipy.spatial
import matplotlib
import matplotlib.pyplot as plt
#####
##### Data Generating Functions #####
#####
VAR_MEASUREMENT_NOISE = 1

def generate_sensors(k = 7, d = 2):
    """
    Generate sensor locations.
    Input:
    k: The number of sensors.
    d: The spatial dimension.
    Output:
    sensor_loc: k * d numpy array.
    """
    sensor_loc = 100*np.random.randn(k,d)
    return sensor_loc

def generate_data(sensor_loc, k = 7, d = 2,
                  n = 1, original_dist = True, sigma_s = 100):
    """
    Generate the locations of n points and distance measurements.

    Input:
    sensor_loc: k * d numpy array. Location of sensor.
    k: The number of sensors.
    d: The spatial dimension.
    n: The number of points.
    original_dist: Whether the data are generated from the original distribution.
    sigma_s: the standard deviation of the distribution that generate each object location.

    Output:
    obj_loc: n * d numpy array. The location of the n objects.
    distance: n * k numpy array. The distance between object and the k sensors.
    """
    assert k, d == sensor_loc.shape

    obj_loc = sigma_s*np.random.randn(n, d)
    if not original_dist:
        obj_loc = sigma_s*np.random.randn(n, d)+([300,300])

    distance = scipy.spatial.distance.cdist(obj_loc,
                                            sensor_loc,
                                            metric='euclidean')
    distance += np.random.randn(n, k)
    return obj_loc, distance

def generate_data_given_location(sensor_loc, obj_loc, k = 7, d = 2):
    """
    Generate the distance measurements given location of a single object and sensor
    .

    Input:
    obj_loc: 1 * d numpy array. Location of object
    sensor_loc: k * d numpy array. Location of sensor.
    k: The number of sensors.
    d: The spatial dimension.
    """
    pass
```

4. b)

```
In [9]: ##### Part b #####
#####
##### Gradient Computing and MLE #####
#####
##### def compute_gradient_of_likelihood(single_obj_loc, sensor_loc,
#####                                     single_distance):
#####
##### Compute the gradient of the loglikelihood function for part a.
#####
##### Input:
##### single_obj_loc: 1 * d numpy array.
##### Location of the single object.
#####
##### sensor_loc: k * d numpy array.
##### Location of sensor.
#####
##### single_distance: k dimensional numpy array.
##### Observed distance of the object.
#####
##### Output:
##### grad: d-dimensional numpy array.
#####
#####
#     grad = np.zeros_like(single_obj_loc)
#     #Your code: implement the gradient of loglikelihood
#     d = len(single_obj_loc[0])
#     k = len(single_distance)
#
#     for i in range(d):
#         #grad[0][i] = np.sum([2*(1 - single_distance[j]*(1/np.linalg.norm(sensor_loc[j]-single_obj_loc[0]))) * (sensor_loc[j][i] - single_obj_loc[0][i]) for j in range(k)])
#         for j in range(k):
#             norm = np.linalg.norm(sensor_loc[j]-single_obj_loc[0]) #np.sqrt(np.square(sensor_loc[j]-single_obj_loc[0]))
#             scalar = 2*(sensor_loc[j][i] - single_obj_loc[0][i])
#
#             grad[0][i] += scalar*(norm - single_distance[j])/norm
#     grad = np.zeros_like(single_obj_loc)
#
#
#     for i in range(len(single_distance)):
#         actual_dist = np.sqrt(np.sum((sensor_loc[i] - single_obj_loc) ** 2))
#         observed_dist = single_distance[i]
#         grad_i = single_obj_loc - sensor_loc[i]
#         grad_i *= ((actual_dist - observed_dist) / actual_dist)
#         grad += grad_i
#
#     grad *= -1
#
#     return grad
#
#def find_mle_by_grad_descent_part_b(initial_obj_loc,
#                                     sensor_loc, single_distance, lr=0.001, num_iters = 10000):
#    """
#    Compute the gradient of the loglikelihood function for part a.
#
#    Input:
#    initial_obj_loc: 1 * d numpy array.
#    Initialized Location of the single object.
```

```
In [10]: ##### MAIN #####
##### Your code: set some appropriate learning rate here #####
def run(lr):
    print("Learning Rate is:", lr)
    np.random.seed(0)
    sensor_loc = generate_sensors()
    obj_loc, distance = generate_data(sensor_loc)
    single_distance = distance[0]
    print('The real object location is')
    print(obj_loc)
    # Initialized as [0,0]
    initial_obj_loc = np.array([[0.,0.]])
    estimated_obj_loc = find_mle_by_grad_descent_part_b(initial_obj_loc,
        sensor_loc, single_distance, lr=lr, num_iters = 10000)
    print('The estimated object location with zero initialization is')
    print(estimated_obj_loc)

    # Random initialization.
    initial_obj_loc = np.random.randn(1,2)
    estimated_obj_loc = find_mle_by_grad_descent_part_b(initial_obj_loc,
        sensor_loc, single_distance, lr=lr, num_iters = 10000)
    print('The estimated object location with random initialization is')
    print(estimated_obj_loc)
    print('\n\n')

lrs = [1.0, 0.01, 0.001, 0.0001]
for lr in lrs:
    run(lr)
```

```
Learning Rate is: 1.0
The real object location is
[[ 44.38632327  33.36743274]]

/opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:43: RuntimeWarning:
overflow encountered in square
/opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:46: RuntimeWarning:
invalid value encountered in double_scalars
```

```
The estimated object location with zero initialization is
[[ nan  nan]]
The estimated object location with random initialization is
[[ nan  nan]]
```

```
Learning Rate is: 0.01
The real object location is
[[ 44.38632327  33.36743274]]

The estimated object location with zero initialization is
[[ 43.07188433  32.71217817]]
The estimated object location with random initialization is
[[ 43.07188433  32.71217817]]
```

```
Learning Rate is: 0.001
The real object location is
[[ 44.38632327  33.36743274]]

The estimated object location with zero initialization is
[[ 43.07188426  32.71217807]]
The estimated object location with random initialization is
[[ 43.07188426  32.71217807]]
```

```
Learning Rate is: 0.0001
The real object location is
[[ 44.38632327  33.36743274]]

The estimated object location with zero initialization is
[[ 39.61604031  27.77608078]]
The estimated object location with random initialization is
[[ 39.67888838  27.90249244]]
```

4. c)

```
In [11]: ##### Part c #####
##### def log_likelihood(obj_loc, sensor_loc, distance):
"""
    This function computes the log likelihood (as expressed in Part a).
    Input:
        obj_loc: shape [1,2]
        sensor_loc: shape [7,2]
        distance: shape [7]
    Output:
        The log likelihood function value.
"""

# Your code: compute the log likelihood
#     func_value = -1*np.sum(
#         [np.square(np.linalg.norm(sensor_loc[i] - obj_loc[0]) - distance[i])
#          for i in range(len(distance))])

#     return func_value
func_value = 0.0

for i in range(len(distance)):
    actual_dist = np.sqrt(np.sum((sensor_loc[i] - obj_loc) ** 2))
    observed_dist = distance[i]
    func_value += (actual_dist - observed_dist) ** 2

func_value *= -1

return func_value
```

```
In [12]: #####
##### Compute the function value at local minimum for all experiments.#####
##### Run 4 #####
def run4(num_sensors):
    num_sensors = 7

    np.random.seed(100)
    sensor_loc = generate_sensors(k=num_sensors)

    # num_data_replicates = 10
    num_gd_replicates = 100

    obj_locs = [[[i,i]] for i in np.arange(0,1000,100)]

    func_values = np.zeros((len(obj_locs),10, num_gd_replicates))
    # record sensor_loc, obj_loc, 100 found minimas
    minimas = np.zeros((len(obj_locs), 10, num_gd_replicates, 2))
    true_object_locs = np.zeros((len(obj_locs), 10, 2))

    for i, obj_loc in enumerate(obj_locs):
        for j in range(10):
            obj_loc, distance = generate_data_given_location(sensor_loc, obj_loc,
                                                               k = num_sensors, d = 2)
            true_object_locs[i, j, :] = np.array(obj_loc)

            for gd_replicate in range(num_gd_replicates):
                initial_obj_loc = np.random.randn(1,2)* (100 * i+1)
                obj_loc = find_mle_by_grad_descent_part_b(initial_obj_loc,
                                                sensor_loc, distance[0], lr=0.1, num_iters = 1000)
                minimas[i, j, gd_replicate, :] = np.array(obj_loc)
                func_value = log_likelihood(obj_loc, sensor_loc, distance[0])
                func_values[i, j, gd_replicate] = func_value

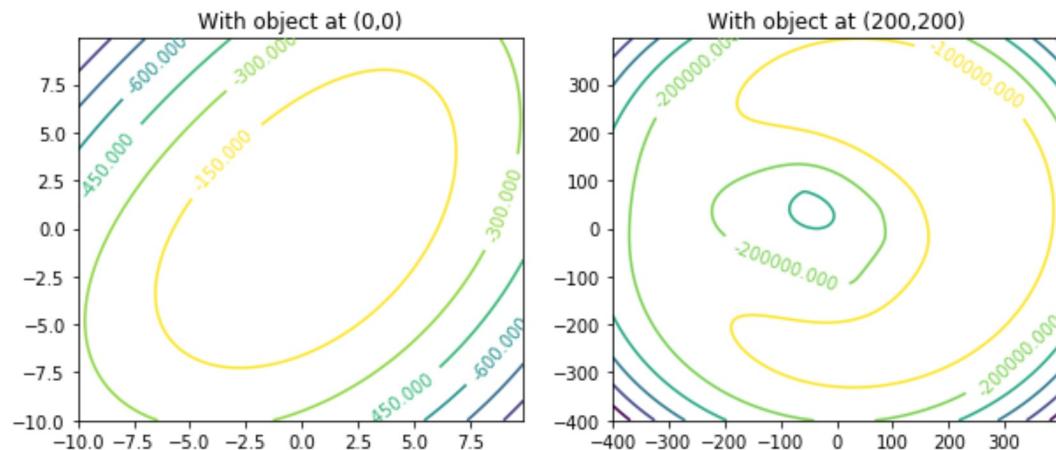
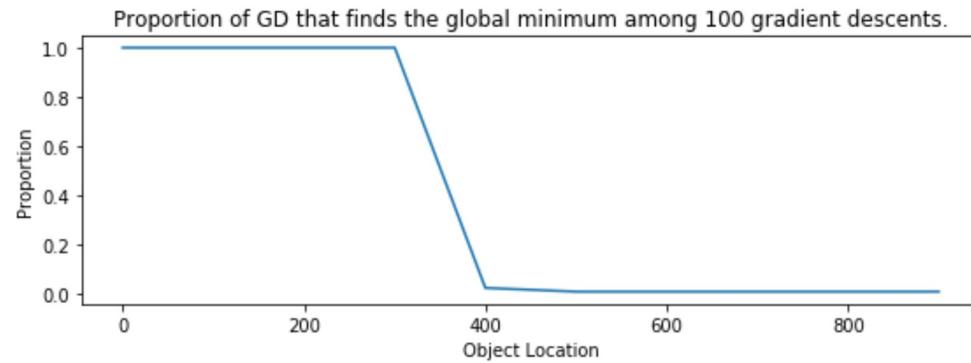
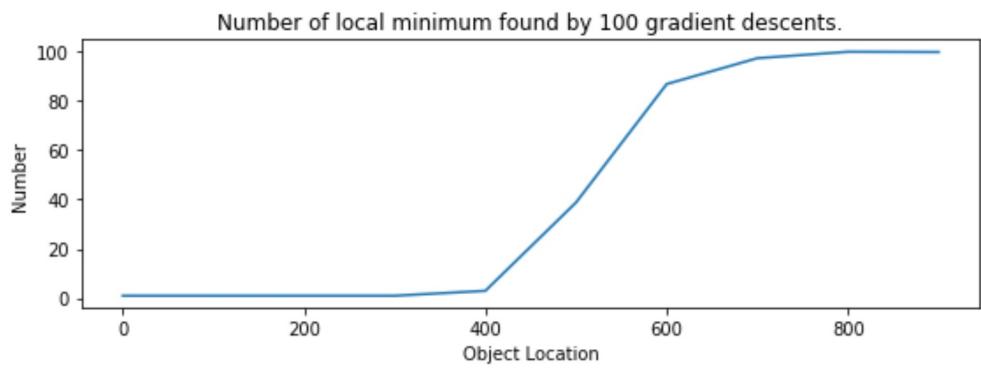
    #####
    ##### Calculate the things to be plotted. #####
    #####
    local_mins = [[np.unique(func_values[i,j]).round(decimals=2)) for j in range(10)]
    for i in range(10)]
    num_local_min = [[len(local_mins[i][j]) for j in range(10)] for i in range(10)]
    proportion_global = [[sum(func_values[i,j].round(decimals=2) == min(local_mins[i][j]))*1.0/100 \
                           for j in range(10)] for i in range(10)]

    num_local_min = np.array(num_local_min)
    num_local_min = np.mean(num_local_min, axis = 1)

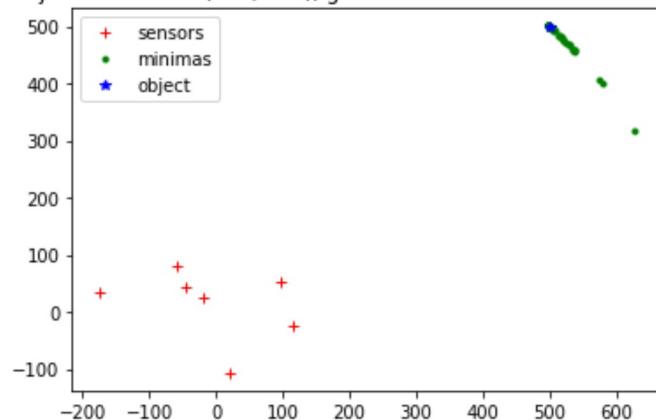
    proportion_global = np.array(proportion_global)
    proportion_global = np.mean(proportion_global, axis = 1)

    #####
    ##### Plots. #####
    #####
    fig, axes = plt.subplots(figsize=(8, 6), nrows=2, ncols=1)
    fig.tight_layout()
    plt.subplot(211)

    plt.plot(np.arange(0,1000,100), num_local_min)
    plt.title('Number of local minimum found by 100 gradient descents.')
    plt.xlabel('Object Location')
    plt.ylabel('Number')
    plt.savefig('num_obj.png')
    # Proportion of gradient descents that find the local minimum of minimum value.
```

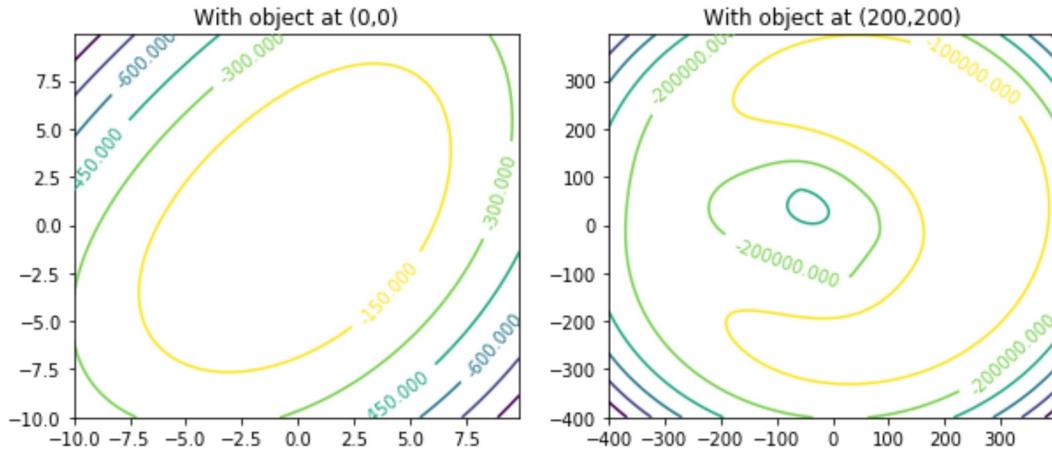
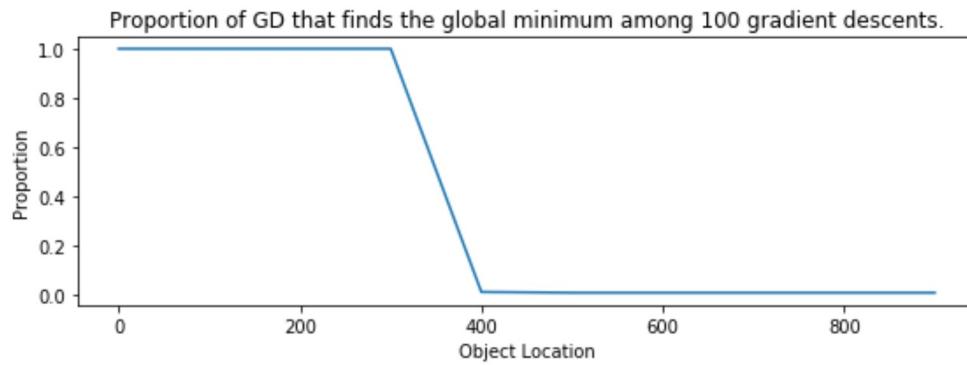
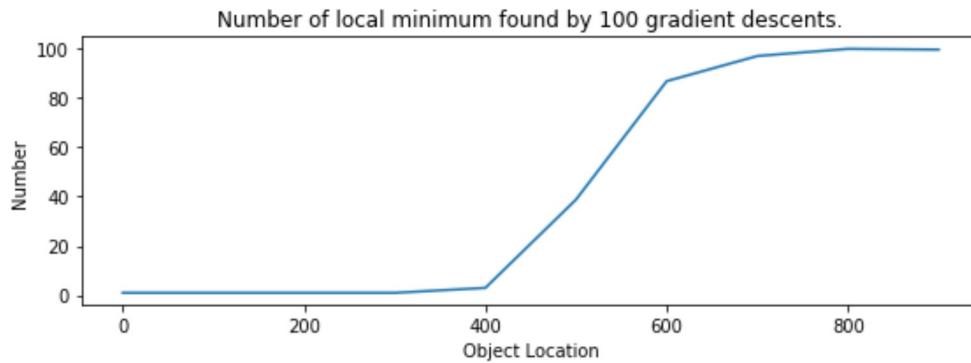


object at location (500, 500), gradient descent recovered locations

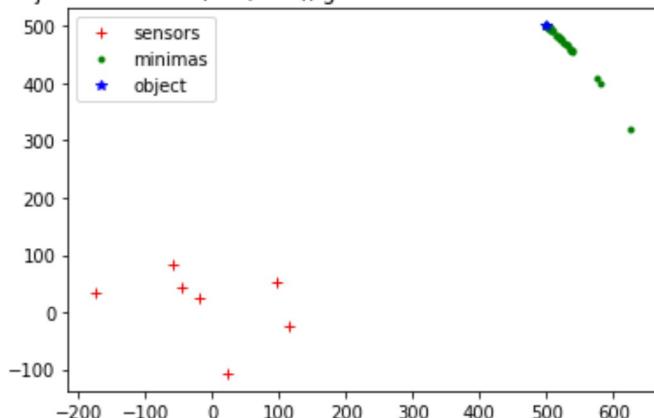


4. d)

```
In [14]: VAR_MEASUREMENT_NOISE = 0.01 # Used by generate_data_given_location function
run4(num_sensors=7)
```



object at location (500, 500), gradient descent recovered locations



4. e)

```
In [15]: VAR_MEASUREMENT_NOISE = 1 # Used by generate_data_given_location function
run4(num_sensors=20)

-----
KeyboardInterrupt                                     Traceback (most recent call last)
<ipython-input-15-e08484daed3e> in <module>()
      1 VAR_MEASUREMENT_NOISE = 1 # Used by generate_data_given_location function
--> 2 run4(num_sensors=20)

<ipython-input-12-8d0fd7e8f169> in run4(num_sensors)
    27     initial_obj_loc = np.random.randn(1,2)* (100 * i+1)
    28     obj_loc = find_mle_by_grad_descent_part_b(initial_obj_loc,
c,
--> 29             sensor_loc, distance[0], lr=0.1, num_iters = 10
00)
    30     minimas[i, j, gd_replicate, :] = np.array(obj_loc)
    31     func_value = log_likelihood(obj_loc, sensor_loc, distanc
e[0])

<ipython-input-9-a5a12d7e7dcc> in find_mle_by_grad_descent_part_b(initial_obj_loc, sensor_loc, single_distance, lr, num_iters)
    72     # Your code: do gradient descent
    73     for i in range(num_iters):
--> 74         grad = compute_gradient_of_likelihood(obj_loc, sensor_loc, singl
e_distance)
    75         obj_loc = obj_loc + grad*lr
    76

<ipython-input-9-a5a12d7e7dcc> in compute_gradient_of_likelihood(single_obj_loc, sensor_loc, single_distance)
    41
    42     for i in range(len(single_distance)):
--> 43         actual_dist = np.sqrt(np.sum((sensor_loc[i] - single_obj_loc) **
2))
    44         observed_dist = single_distance[i]
    45         grad_i = single_obj_loc - sensor_loc[i]

/opt/conda/lib/python3.6/site-packages/numpy/core/fromnumeric.py in sum(a, axis, dtype, out, keepdims)
    1812         return sum(axis=axis, dtype=dtype, out=out, **kwargs)
    1813     return _methods._sum(a, axis=axis, dtype=dtype,
-> 1814                         out=out, **kwargs)
    1815
    1816

KeyboardInterrupt:
```

4. f)

In [16]:

```
#####
##### Gradient Computing and MLE #####
#####
##### def compute_grad_likelihood(sensor_loc, obj_loc, distance):
#####
# Compute the gradient of the loglikelihood function for part f.

Input:
sensor_loc: k * d numpy array.
Location of sensors.

obj_loc: n * d numpy array.
Location of the objects.

distance: n * k dimensional numpy array.
Observed distance of the object.

Output:
grad: k * d numpy array.
#####
grad = np.zeros(sensor_loc.shape)
# Your code: finish the grad loglike
for i, sens in enumerate(sensor_loc):
    grad[i] = compute_gradient_of_likelihood(sens.reshape((1,-1)), obj_loc, distance[:, i])

return grad

def find_mle_by_grad_descent(initial_sensor_loc,
                             obj_loc, distance, lr=0.001, num_iters = 1000):
#####
# Compute the gradient of the loglikelihood function for part f.

Input:
initial_sensor_loc: k * d numpy array.
Initialized Location of the sensors.

obj_loc: n * d numpy array. Location of the n objects.

distance: n * k dimensional numpy array.
Observed distance of the n object.

Output:
sensor_loc: k * d numpy array. The mle for the location of the object.

#####
sensor_loc = initial_sensor_loc
# Your code: finish the gradient descent
for i in range(num_iters):
    grad = compute_grad_likelihood(sensor_loc, obj_loc, distance)
    sensor_loc += lr*grad

return sensor_loc
#####
##### Gradient Computing and MLE #####
#####
np.random.seed(0)
sensor_loc = generate_sensors()
obj_loc, distance = generate_data(sensor_loc, n = 100)
print('The real sensor locations are')
print(sensor_loc)
# Initialized as zeros.
initial_sensor_loc = np.zeros((7,2)) #np.random.randn(7,2)
```

The real sensor locations are
[[176.4052346 40.01572084]
[97.87379841 224.08931992]
[186.75579901 -97.72778799]
[95.00884175 -15.13572083]
[-10.32188518 41.05985019]
[14.40435712 145.4273507]
[76.10377251 12.16750165]]

The predicted sensor locations are
[[176.28223286 39.76550241]
[97.73185804 224.12127837]
[186.68461116 -97.5920796]
[94.90639305 -15.16741152]
[-10.17511568 40.93691838]
[14.19497946 145.42410633]
[76.11863016 12.24807699]]

The MSE for Case 1 is 7.863128736936463

The MSE for Case 2 is 7.227391044803935

The MSE for Case 2 (if we knew mu is [300,300]) is 6.914124924780527

Comments for Question 4

4. B

Learning rate of 1 is unreasonable. The best learning rate was 0.01 for both random and zero initialization.

4. C

Do you find any patterns?

As our optimal point is moved farther away from the sensors our gradient descent to get to the optimal position becomes worse. The insight we get from slight differences between our sensors which we use in GD to determine optimal direction become useless.

4. D

Do you find any patterns?

“Reduce variance of the measurement noise”

The difference did not amount to much. This is due to the fact that the variance which we have across all of our measurement noise only scales our gradient descent. It is worth noting that our algorithm may have to take smaller steps to guarantee converge due to this scaling potential.

4. E

Do you find any patterns?

More sensors means that it is easier to locate our object. Having a greater number of sensors means that our gradient descent is better able to distinguish an optimal direction to find the global minimum.