

Q1

$$\textcircled{1} \quad f(x+p) = f(x) + \nabla f(x)^T p + O(\|p\|^2)$$

From Taylor's theorem, the k^{th} order Taylor polynomial as:

$$f(b) = f(a) + f'(a)(b-a) + \frac{f''(a)}{2!}(b-a)^2 + \cdots + \frac{f^{(k)}(a)}{k!}(b-a)^k \\ + h_k(b)(b-a)^k, \lim_{b \rightarrow a} h_k(b) = 0$$

Therefore we can get :

The best linear approximation for $f(b)$ is:

$$f(b) = f(a) + f'(a)(b-a) + R, R \text{ represents the error term between approximation to real } f(b) \text{ value.}$$

For a multivariable version,

$$f(b) = f(b_1, b_2, b_3, \dots, b_n)$$

then the linear approximation (first order) of $f(b)$ is:

$$f(b) = f(a) + \nabla f(a)^T (b-a) + R$$

where $\nabla f(a)$ represents a $1 \times n$ horizontal vector of partial derivatives of $f(a_1, a_2, \dots, a_n)$

substitute b with $x+p$, a with x , we get:

$$f(x+p) = f(x) + \nabla f(x)^T p + R$$

Since $R = \frac{f''(x)}{2!}(p)^2 + \cdots + \frac{f^{(k)}(x)}{k!} p^k$ and $\lim_{x+p \rightarrow x} R = 0$,

$R = O(\|p\|^2)$ (dominated by quadratic term),

$$f(x+p) = f(x) + \nabla f(x)^T p + O(\|p\|^2)$$

$$\textcircled{2} f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + O(\|p\|^3)$$

Similar to \textcircled{1}, with second order Taylor polynomial
or

$$f(b) = f(a) + f'(a)(b-a) + \frac{f''(a)}{2}(b-a)^2 + R, R \text{ represents}$$

the error term between approximation to real $f(b)$ value

For a multivariable version, $f(b) = f(b_1, b_2, b_3, \dots, b_n)$

From above, since $f(b)$ is scalar, first derivative $f'(a)$ is a $1 \times n$ matrix denotes by $\nabla f(a)$,
for the second derivative of $f(a)$, we can take
the matrix of partial derivative of the vector $\nabla f(a)$
which written as $\nabla^2 f(a)$, is an $n \times n$ matrix

Therefore, we can rewrite a single variable quadratic expression $\frac{1}{2}(b-a)^T f''(a)(b-a)$ into multivariable version $\frac{1}{2}(b-a)^T \nabla^2 f(a)(b-a)$

Combine the linear expression we derived from part \textcircled{1},
substitute b with $x+p$, a with x , we get:

$$f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + R$$

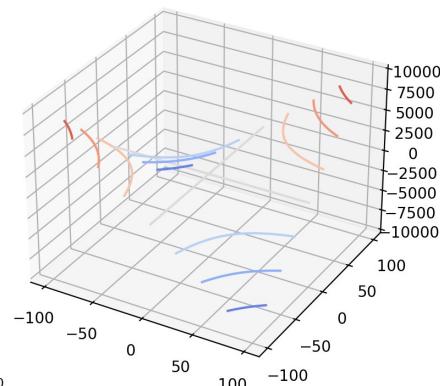
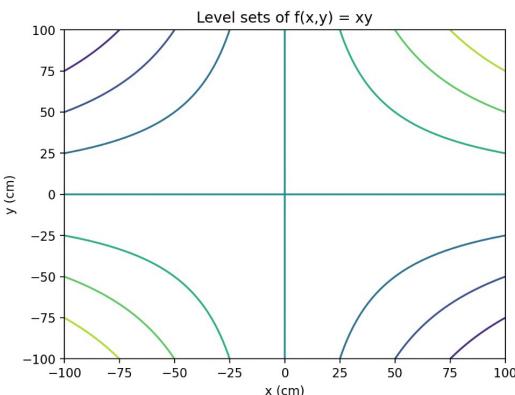
Since $R = \frac{f''(x)}{3!}(p)^3 + \dots + \frac{f^k(x)}{k!} p^k$ and $\lim_{x+p \rightarrow x} R = 0$,

$R = O(\|p\|^3)$ (dominated by cubic term),

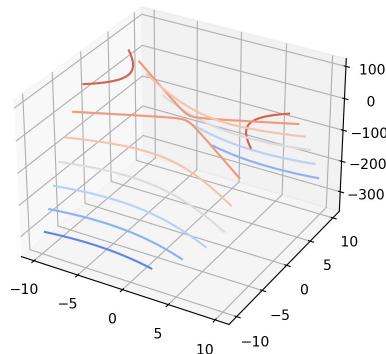
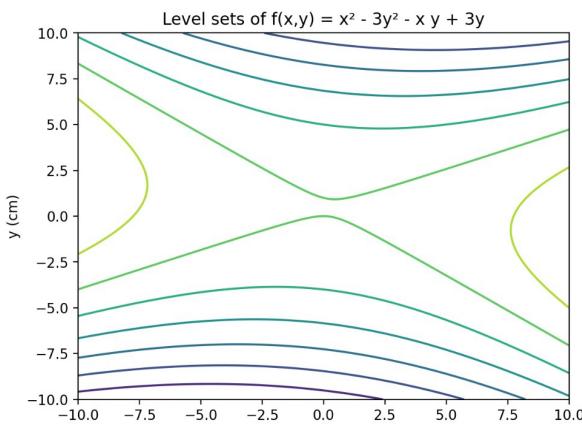
$$f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + O(\|p\|^3)$$

Q₂

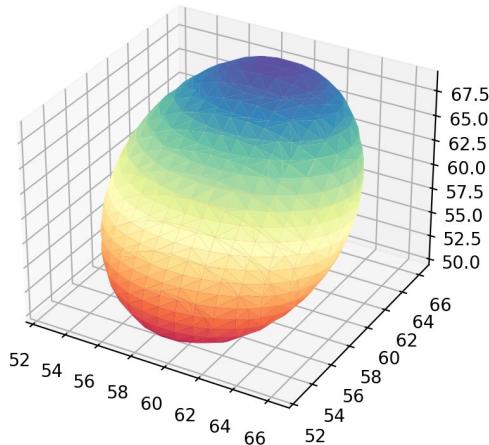
(1)



(2)



(3)



Q3 Convex function $f(x)$, by definition that on an interval $[a, b]$ if for any two points x_1 and x_2 in $[a, b]$ and with λ $0 < \lambda < 1$:

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \quad ①$$

Let x^* be a local minima at convex function $f(x)$, which means $f(x^*) \geq f(x)$ for every x^* in $B(x, \delta) [\delta > 0]$.

Suppose there is a point y , $f(y) < f(x)$

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

for $0 < \lambda < 1$, if we pick a large enough λ , $\lambda x + (1-\lambda)y$ will fall in the interval of $(x-\delta, x+\delta)$, we have:

$$\lambda f(x) + (1-\lambda)f(y) \geq f(\lambda x + (1-\lambda)y) \geq f(x)$$

$$\lambda f(x) + (1-\lambda)f(y) \geq f(x)$$

$f(y) \geq f(x)$, so $f(x)$ is global minima

Q4

Given $f(x)$ being positive, since x^* maximize $f(x)$,
 $f(x^*) \geq f(x) > 0$, x being any input for $f(x)$.
since $-\log x$ monotonically decrease in $(0, +\infty)$,
 $-\log(f(x^*)) \leq f(x)$, which means x^*
minimize $-\log(f(x))$.

if a, b on $f(x)$, b being any point other than
 a , and $-\log(a) \leq -\log(b)$, given $f(x)$ is positive,
 a must be the global maxima on $f(x)$.
Given a is achieved by an input x^* to $f(x)$,
therefore any input x^* that minimize $-\log(f(x))$
also maximize $f(x)$

$$② L(\mu, \Sigma) = \prod_{i=1}^k N(x_i | \mu, \Sigma)$$

$$= \prod_{i=1}^k (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right)$$

$$= 2\pi^{-\frac{k n}{2}} |\Sigma|^{-\frac{k}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^k (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right)$$

$$-\ln(L(\mu, \Sigma)) = \ln 2\pi^{-\frac{k n}{2}} - \ln |\Sigma|^{-\frac{k}{2}} - \ln \exp\left(-\frac{1}{2} \sum_{i=1}^k (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right)$$

$$= \frac{k}{2} + \frac{1}{2} \sum_{i=1}^k (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) + C$$

$$-\frac{\partial \ln L}{\partial \mu} = \frac{\partial}{\partial \mu} \left(\frac{1}{2} \sum_{i=1}^k (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$(x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \overbrace{(x_i^T - \mu^T)}^{\leftarrow} \Sigma^{-1} (x_i - \mu)$$

$$= x^T \Sigma^{-1} x - \underbrace{x^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x}_{\cancel{x^T \Sigma^{-1} \mu}} + \mu^T \Sigma^{-1} \mu$$

$$\cancel{x^T \Sigma^{-1} \mu}$$

$$\sum_{i=1}^k (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \sum_{i=1}^k x_i^T \Sigma^{-1} x_i - \cancel{2 \sum_{i=1}^k \mu^T \Sigma^{-1} x_i} + k \mu^T \Sigma^{-1} \mu$$

$$-\frac{\partial \ln L}{\partial \mu} = \frac{\partial}{\partial \mu} \frac{1}{2} \left(\sum_{i=1}^k x_i^\top \Sigma^{-1} x_i - 2 \sum_{i=1}^k \mu^\top \Sigma^{-1} x_i + k \mu^\top \Sigma^{-1} \mu \right)$$

$$= \frac{\partial}{\partial \mu} \left(- \sum_{i=1}^k x_i^\top \Sigma^{-1} \mu \right) + \frac{1}{2} \frac{\partial}{\partial \mu} k \mu^\top \Sigma^{-1} \mu$$

$$\textcircled{1} = - \sum_{i=1}^k \Sigma^{-1} x_i \quad [f(x) = A^\top x, \frac{df}{dx} = A]$$

$$\Sigma^{-1 T} = \Sigma^{-1} \Rightarrow \textcircled{1} = - \sum_{i=1}^k \Sigma^{-1} x_i^\top = - \sum_{i=1}^k \Sigma^{-1} x_i$$

$$\textcircled{2} = k \sum_i \mu \quad [f(x) = x^\top A x, \frac{df}{dx} = Ax + A^\top x = 2Ax \quad (A = A^\top)]$$

$$-\frac{\partial \ln L}{\partial \mu} = - \sum_{i=1}^k \Sigma^{-1} x_i^\top + k \sum_i \mu$$

$$-\frac{\partial \ln L}{\partial \Sigma} = \underbrace{\frac{k}{2} \frac{\partial}{\partial \Sigma} \ln |\Sigma|}_{\textcircled{1}} + \underbrace{\frac{1}{2} \frac{\partial}{\partial \Sigma} \sum_{i=1}^k (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}_{\textcircled{2}}$$

$$\textcircled{1} = \frac{k}{2} \Sigma^{-1}$$

$$\textcircled{2} = \frac{1}{2} \sum_{i=1}^k (-\Sigma^{-1} (x_i - \mu) (x_i - \mu)^T \Sigma^{-1})$$

$$= \frac{k}{2} \Sigma^{-1} + \frac{1}{2} \sum_{i=1}^k (-\Sigma^{-1} (x_i - \mu) (x_i - \mu)^T \Sigma^{-1})$$

$$-\frac{\partial \ln L}{\partial \Sigma} = \frac{1}{2} \Sigma^{-1} \left(k\Sigma - \sum_{i=1}^k (x_i - \mu) (x_i - \mu)^T \right) \Sigma^{-1}$$

(3)

$$\nabla - \ln L(\mu) = \frac{\partial \ln L}{\partial \mu} = -\sum^{-1} \sum_{i=1}^k x^i + k \sum^{-1} \mu$$

$$\nabla - \ln L(\Sigma) = \frac{\partial \ln L}{\partial \Sigma} = \frac{1}{2} \sum^{-1} \left(k \sum - \sum_{i=1}^k (x^i - \mu)(x^i - \mu)^T \right) \sum^{-1}$$

$$-\nabla \ln L(\mu) = 0$$

$$\Rightarrow \frac{\sum^{-1} \sum_{i=1}^k x^i}{k \sum^{-1} \mu} = 0 \quad \mu = \frac{1}{k} \sum_{i=1}^k x^i = \bar{x}$$

$$-\nabla \ln L(\Sigma) = 0$$

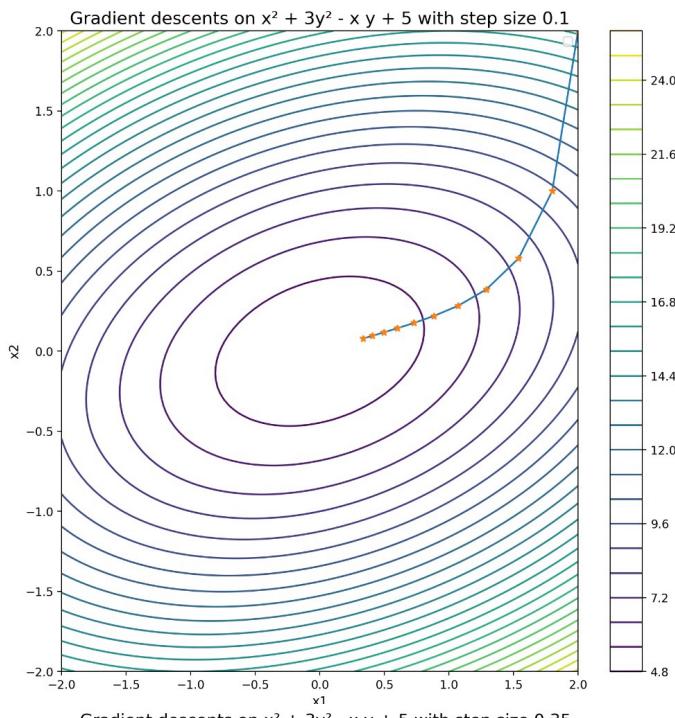
$$\Rightarrow \left(k \sum - \sum_{i=1}^k (x - \mu)(x - \mu)^T \right) = 0$$

$$\sum = \frac{1}{k} \sum_{i=1}^k (x - \mu)(x - \mu)^T$$

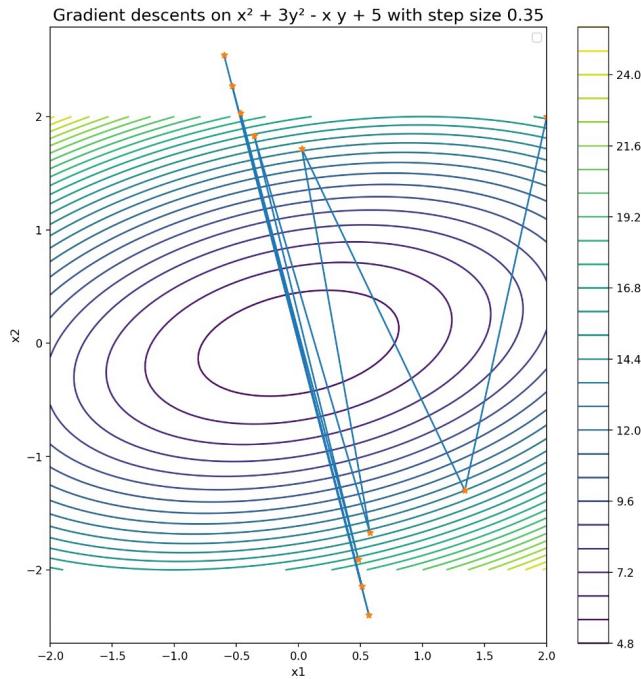
(4) Since $L(\theta)$ is convex function and
 θ is the point that locally maximized
 $L(\theta)$ by doing gradient descent, so
 θ globally maximizes $L(\theta)$

Q5

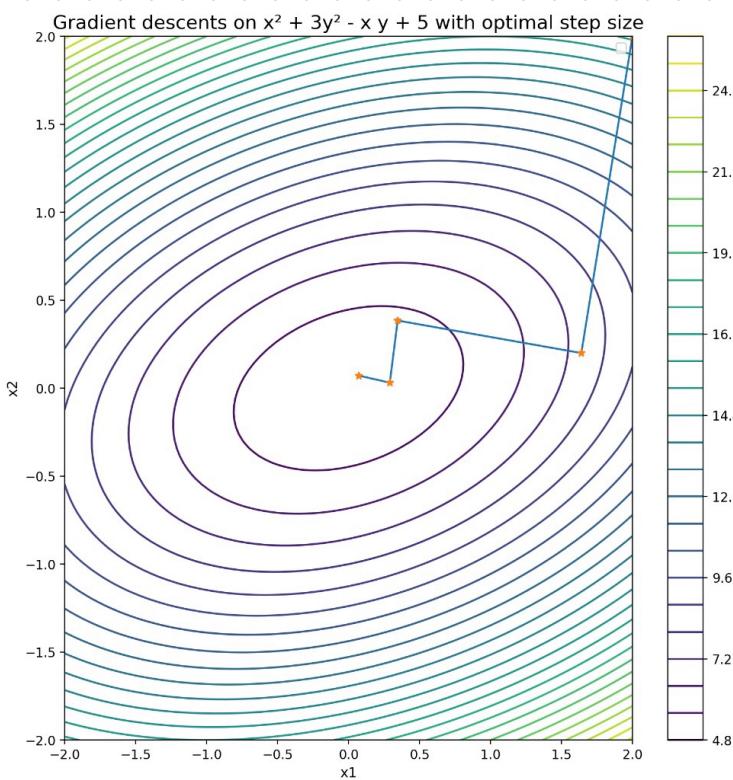
(1)



(2)



(3)



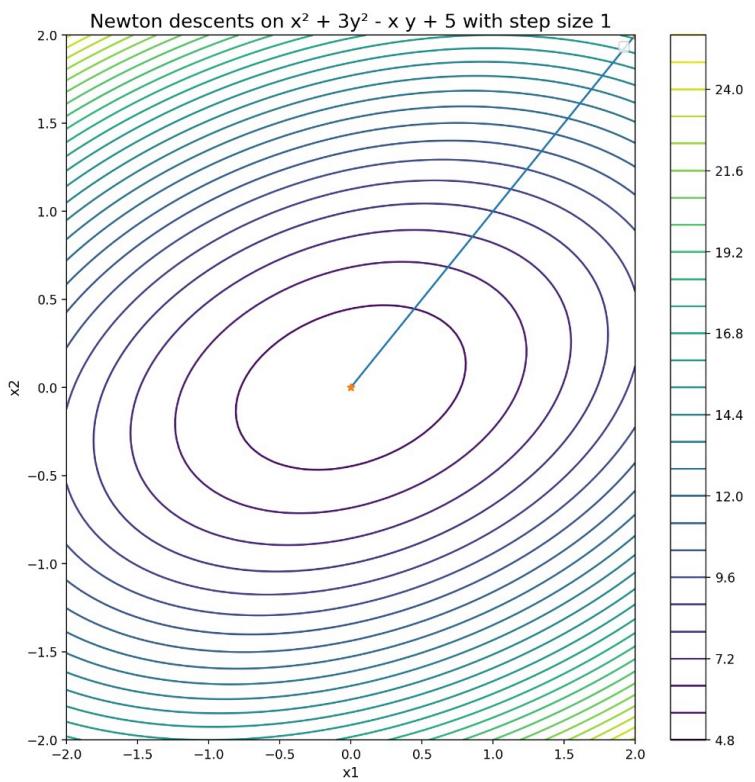
To find optimal step size, I loop through all step value between 0.01 to 10, find the step value α such:

$$\boxed{\nabla f(x^{(k)} + p^{(k)} \cdot \alpha)^T \cdot p^{(k)}} \text{ closest to } 0$$

$$p^{(k+1)}$$

which means the current step direction and next step direction is closest to orthogonal

(4)



Q6

symmetric

① if $Q \in \mathbb{R}^n$ and Q is a positive semi-definite matrix, then eigenvalues of Q are real and positive

Proof of Q 's eigenvalues are real:

if Q is a positive, semi-definite matrix (Also as hermitian matrix), then

$\boxed{Q^T = Q}$, $\boxed{Qx = \lambda x}$, λ is an eigenvalue of Q , x is a eigenvector to eigenvalue λ

$$\Rightarrow \bar{x}^T Q x = \bar{x}^T \lambda x = \bar{\lambda} \bar{x}^T x = \bar{\lambda} \|x\|^2 \quad ①$$

$$\text{Also, } x^T Q x = (Qx)^T \bar{x} = \bar{x}^T Q^T \bar{x} \quad ②$$

$$\Rightarrow x^T Q^T x = \lambda \|x\|^2 \quad ③$$

with complete conjugate on ③, we have

$$x^T Q^T \bar{x} = \bar{\lambda} \|x\|^2$$

$$\text{Since } Q^T = Q$$

$$x^T Q \bar{x} = \bar{\lambda} \|x\|^2 = x^T \bar{\lambda} x = \bar{\lambda} \|x\|^2, \bar{\lambda} = \lambda, \|x\| \neq 0$$

$\therefore \lambda$ is a real number

Proof of λ is positive :

$$Qx = \lambda x \text{ from } ①$$

$$x^T Q x = \lambda x^T x = \lambda \|x\|^2$$

positive because Q is positive semi definite,
 x is also positive bc they are eigenvector

Q. Which means $\lambda \|x\|^2$ λ is positive

since $\|x\|^2$ is positive as well

② From spectrum theorem, a symmetric, positive definite matrix is diagonalizable, which means there exists a matrix D, which

$D = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 \\ \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$ that $\lambda_1 \dots \lambda_n$ are Q's eigenvalues.

Therefore, to pair with these n eigenvalues, there are n eigenvectors.

arbitrary

Let v_1, v_2 be the two of N eigenvectors with eigen value λ_1, λ_2

$$Qv_1 = \lambda_1 v_1, Qv_2 = \lambda_2 v_2, Q = Q^T$$

$$Qv_1 \cdot v_2 = v_1 \cdot Q^T v_2 = v_1 \cdot Q v_2$$

$$\lambda_1 v_1 \cdot v_2 = v_1 \cdot \lambda_2 v_2$$

$$(\lambda_1 - \lambda_2) v_1 \cdot v_2 = 0, \lambda_1 \neq \lambda_2$$

$\Rightarrow v_1 \cdot v_2 = 0 \Rightarrow v_1$ and v_2 are orthogonal to each other