

CSE 512 Machine Learning

HW #3

Enbo Yu

113094714

12 Oct (3 days extension)

1. Gradient properties.

(a) Linearity. If $h(x) = \alpha f(x) + \beta g(x)$, then $\nabla h(x) = \alpha \nabla f(x) + \beta \nabla g(x)$.

$$h(x) = \alpha f(x) + \beta g(x),$$

We know that $\nabla f(x)$ is the gradient of f at x .

so for $\nabla h(x)$:

$$\nabla h(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} h(x_1) \\ \frac{\partial}{\partial x_2} h(x_2) \\ \dots \dots \\ \frac{\partial}{\partial x_n} h(x_n) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} (\alpha f(x_1) + \beta g(x_1)) \\ \frac{\partial}{\partial x_2} (\alpha f(x_2) + \beta g(x_2)) \\ \dots \dots \\ \frac{\partial}{\partial x_n} (\alpha f(x_n) + \beta g(x_n)) \end{bmatrix}$$

by the definition of Linearity of differentiation,

$$\frac{d}{dx} (\alpha \cdot f(x) + \beta \cdot g(x)) = \alpha \cdot f'(x) + \beta \cdot g'(x).$$

$$\text{so } \nabla h(x) = \begin{bmatrix} \alpha \cdot \frac{\partial}{\partial x_1} f(x_1) + \beta \frac{\partial}{\partial x_1} g(x_1) \\ \alpha \cdot \frac{\partial}{\partial x_2} f(x_2) + \beta \frac{\partial}{\partial x_2} g(x_2) \\ \dots \dots \\ \alpha \frac{\partial}{\partial x_n} f(x_n) + \beta \frac{\partial}{\partial x_n} g(x_n) \end{bmatrix} = \begin{bmatrix} \alpha \cdot \frac{\partial}{\partial x_1} f(x_1) \\ \alpha \cdot \frac{\partial}{\partial x_2} f(x_2) \\ \dots \\ \alpha \frac{\partial}{\partial x_n} f(x_n) \end{bmatrix} + \begin{bmatrix} \beta \frac{\partial}{\partial x_1} g(x_1) \\ \beta \frac{\partial}{\partial x_2} g(x_2) \\ \dots \\ \beta \frac{\partial}{\partial x_n} g(x_n) \end{bmatrix}$$

$$= \alpha \cdot \begin{bmatrix} \frac{\partial}{\partial x_1} f(x_1) \\ \frac{\partial}{\partial x_2} f(x_2) \\ \dots \\ \frac{\partial}{\partial x_n} f(x_n) \end{bmatrix} + \beta \cdot \begin{bmatrix} \frac{\partial}{\partial x_1} g(x_1) \\ \frac{\partial}{\partial x_2} g(x_2) \\ \dots \\ \frac{\partial}{\partial x_n} g(x_n) \end{bmatrix} = \alpha \cdot \nabla f(x) + \beta \cdot \nabla g(x)$$

$$\text{so } \nabla h(x) = \alpha \nabla f(x) + \beta \nabla g(x).$$

(b) Chain rule. If $g(v) = f(Av)$, then $\nabla g(v) = A^T \nabla f(Av)$.

Based on question (a).

$$\begin{aligned} \nabla g(v) &= \left[\begin{array}{c} \frac{\partial}{\partial v_1} g(v_1) \\ \frac{\partial}{\partial v_2} g(v_2) \\ \vdots \\ \frac{\partial}{\partial v_n} g(v_n) \end{array} \right] \xrightarrow{\text{by chain rule}} \frac{dy}{dx} = f'(g(x)) \cdot g'(x) \left[\begin{array}{c} \frac{\partial}{\partial Av_1} f(Av_1) \cdot \frac{\partial}{\partial v_1}(Av_1) \\ \frac{\partial}{\partial Av_2} f(Av_2) \cdot \frac{\partial}{\partial v_2}(Av_2) \\ \vdots \\ \frac{\partial}{\partial Av_n} f(Av_n) \cdot \frac{\partial}{\partial v_n}(Av_n) \end{array} \right] \\ &= \left[\begin{array}{c} \frac{\partial}{\partial Av_1} f(Av_1) \cdot A_1 \\ \frac{\partial}{\partial Av_2} f(Av_2) A_2 \\ \vdots \\ \frac{\partial}{\partial Av_n} f(Av_n) A_n \end{array} \right] = [A_1, A_2, \dots, A_n] \left[\begin{array}{c} \frac{\partial}{\partial Av_1} f(Av_1) \\ \frac{\partial}{\partial Av_2} f(Av_2) \\ \vdots \\ \frac{\partial}{\partial Av_n} f(Av_n) \end{array} \right] = A^T \nabla f(Av). \end{aligned}$$

so we can get $\nabla g(v) = A^T \nabla f(Av)$.

2. Gradients.

(a) Quadratic function. $f(x) = \frac{1}{2} x^T Q x + p^T x + r$, $Q \in \mathbb{R}^{12 \times 12}$ and Q is symmetric ($Q_{i,j} = Q_{j,i}$).

$$\nabla f(x) = \left[\begin{array}{c} \frac{\partial}{\partial x_1} (\frac{1}{2} x_1^T Q x_1 + p^T x_1 + r) \\ \frac{\partial}{\partial x_2} (\frac{1}{2} x_2^T Q x_2 + p^T x_2 + r) \\ \vdots \\ \frac{\partial}{\partial x_n} (\frac{1}{2} x_n^T Q x_n + p^T x_n + r) \end{array} \right], \text{ by the linearity in Question 1 (a), we can get:}$$

$$\begin{aligned} &= \left[\begin{array}{c} \frac{\partial}{\partial x_1} (\frac{1}{2} x_1^T Q x_1) \\ \frac{\partial}{\partial x_2} (\frac{1}{2} x_2^T Q x_1) \\ \vdots \\ \frac{\partial}{\partial x_n} (\frac{1}{2} x_n^T Q x_1) \end{array} \right] + \left[\begin{array}{c} \frac{\partial}{\partial x_1} (p^T x_1) \\ \frac{\partial}{\partial x_2} (p^T x_1) \\ \vdots \\ \frac{\partial}{\partial x_n} (p^T x_1) \end{array} \right] + \left[\begin{array}{c} \frac{\partial}{\partial x_1} r \\ \frac{\partial}{\partial x_2} r \\ \vdots \\ \frac{\partial}{\partial x_n} r \end{array} \right] = \frac{1}{2} \left[\begin{array}{c} \frac{\partial}{\partial x_1} (x_1^T Q x_1) \\ \frac{\partial}{\partial x_2} (x_2^T Q x_1) \\ \vdots \\ \frac{\partial}{\partial x_n} (x_n^T Q x_1) \end{array} \right] + \left[\begin{array}{c} \frac{\partial}{\partial x_1} p^T x_1 \\ \frac{\partial}{\partial x_2} p^T x_1 \\ \vdots \\ \frac{\partial}{\partial x_n} p^T x_1 \end{array} \right] \\ &\quad \underbrace{\quad}_{\frac{\partial}{\partial x} (\text{constant}) = 0} \quad \downarrow \end{aligned}$$

$Q \in \mathbb{R}^{12 \times 12}$, and $Q_{i,j} = Q_{j,i}$

$$\begin{aligned} Q &= \left[\begin{array}{ccc} Q_{1,1} & \dots & Q_{1,12} \\ \dots & \dots & \dots \\ Q_{12,1} & \dots & Q_{12,12} \end{array} \right], \frac{\partial}{\partial x} x^T Q x = \sum_{i=1}^{12} x_i Q_{i,1} x_j + \sum_{j=1}^{12} x_j Q_{1,j} x_i \\ &\quad \text{so } \nabla f(x) = \frac{1}{2} \cdot (Qx + Qx) + p + 0 \\ &\quad = Qx + p \end{aligned}$$

the dimension is 12.

(b) $f(x) = \frac{1}{\mu} \log(\sum_{i=1}^8 \exp(\mu x_{[i]}))$, $x \in \mathbb{R}^8$. μ is a positive scalar.

let $g(x) = \sum_{i=1}^8 \exp(\mu x_{[i]})$

$$\nabla f(x) = \begin{bmatrix} \frac{\partial}{\partial g(x)} f(x) \cdot \frac{\partial}{\partial x_1} g(x) \\ \frac{\partial}{\partial g(x)} f(x) \cdot \frac{\partial}{\partial x_2} g(x) \\ \dots \\ \frac{\partial}{\partial g(x)} f(x) \cdot \frac{\partial}{\partial x_n} g(x) \end{bmatrix} = \begin{bmatrix} \frac{1}{\mu} \frac{\mu \cdot \exp(\mu x_{[1]})}{\sum_{i=1}^8 \exp(\mu x_{[i]})} \\ \frac{1}{\mu} \frac{\mu \cdot \exp(\mu x_{[2]})}{\sum_{i=1}^8 \exp(\mu x_{[i]})} \\ \dots \\ \frac{1}{\mu} \frac{\mu \cdot \exp(\mu x_{[8]})}{\sum_{i=1}^8 \exp(\mu x_{[i]})} \end{bmatrix} = \begin{bmatrix} \frac{\exp(\mu x_{[1]})}{\sum \exp(\mu x_{[i]})} \\ \frac{\exp(\mu x_{[2]})}{\sum \exp(\mu x_{[i]})} \\ \dots \\ \frac{\exp(\mu x_{[8]})}{\sum \exp(\mu x_{[i]})} \end{bmatrix}$$

we can set the sum as a term $S(x)$,

$$\text{so } \nabla f(x) = \frac{1}{S(x)} \begin{bmatrix} \exp(\mu x) \\ \dots \\ \exp(\mu x_8) \end{bmatrix}, \text{ the dimension is 8.}$$

3. Convex or not convex.

(a) $S = \{x : \sum_i x_i = 0\}$.

Let $m = \{m : \sum_i m_i = 0\} \in S$.

$n = \{n : \sum_i n_i = 0\} \in S$,

set α , $\alpha \in [0, 1]$.

$$\alpha m + (1-\alpha)n = \alpha \sum_i m_i + (1-\alpha) \sum_i n_i = \alpha \cdot 0 + (1-\alpha) \cdot 0 = 0$$

so $\alpha m + (1-\alpha)n \in S$,

so $S = \{x : \sum_i x_i = 0\}$ is convex.

(b) $S = \{(x, y) : x^2 + y^2 = 1\}$

$x^2 + y^2 = 1$ shows the side of a circle whose radius is 1.

so we can assume that

$$m = (x=0, y=1) \text{ and } n = (m=1, n=0)$$

$\alpha \in [0, 1]$, so we can assume $\alpha = 1/3$,

$$\text{so } \alpha m + (1-\alpha)n = \frac{1}{3}(0, 1) + \frac{2}{3}(1, 0)$$

$$= (0, 1/3) + (2/3, 0)$$

$$x^2 + y^2 = (1/3)^2 + (2/3)^2 = 1/9 + 4/9 = 5/9 \neq 1$$

so $\alpha m + (1-\alpha)n \notin S$.

so S is not convex. (Based on the graph of $x^2 + y^2 = 1$),



We also can easily know that S is not convex)

(c) $S = \{x : |x| \leq 1\}$.

let $m = \{m : |m| \leq 1\}$,

$n = \{n : |n| \leq 1\}$.

set α , $\alpha \in [0, 1]$

$$\begin{aligned} |\alpha m + (1-\alpha)n| &\leq |\alpha m| + |(1-\alpha)n| = \alpha|m| + (1-\alpha)|n| \\ &\leq \alpha \cdot 1 + (1-\alpha) \cdot 1 = 1 \end{aligned}$$

so $|\alpha m + (1-\alpha)n| \in S$.

so S is convex.

4. Am I positive semidefinite?

(a) $X = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 4 \\ 3 & 4 & 1 \end{bmatrix}$, X is positive semidefinite for all v , $v^T Hv \geq 0$.

$$\text{so } v^T X v = [v_1 \ v_2 \ v_3] \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 4 \\ 3 & 4 & 1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} = v_1^2 + 2v_1v_2 + 3v_1v_3 + 2v_1v_2 + v_2^2 + 4v_2v_3 + 3v_1v_3 + 4v_2v_3 + v_3^2$$

$$= v_1^2 + v_2^2 + v_3^2 + 4v_1v_2 + 6v_1v_3 + 8v_2v_3$$

we can set $v = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}$, so $v^T = [-1 \ 1 \ -1]$,

$$\text{so } v^T X v = (-1)(-1) + 1(1) + (-1)(-1) - 4 + 6 - 8 = -3 < 0$$

$$\begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}$$

so X is not positive semidefinite, a counter example is $v = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}$

(b) $X = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}$,

$$v^T X v = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} = v_1^2 + 2v_2^2 + 3v_3^2$$

we know that, for any v , $v^2 \geq 0$, so $v_1^2 + 2v_2^2 + 3v_3^2 \geq 0$

so $v^T X v \geq 0$, so X is positive semidefinite.

5. Convex or not convex.

(a) $f(x) = 1/x$, for $x > 0$.

By definition, let $x > 0$, $y > 0$, and $\alpha \in [0, 1]$,

$$f(\alpha x + (1-\alpha)y) = \frac{1}{\alpha x + (1-\alpha)y}$$

$$\alpha f(x) + (1-\alpha)f(y) = \frac{\alpha}{x} + \frac{(1-\alpha)}{y} = \frac{\alpha y + (1-\alpha)x}{xy}$$

We can assume that $f(\alpha x + (1-\alpha)y) > \alpha f(x) + (1-\alpha)f(y)$

$$\text{so } \frac{1}{\alpha x + (1-\alpha)y} > \frac{\alpha y + (1-\alpha)x}{xy} \quad (x>0, y>0)$$

$$xy > (\alpha y + (1-\alpha)x) \cdot (\alpha x + (1-\alpha)y)$$

$$xy > \alpha^2 xy + \alpha y \cdot (1-\alpha)y + (1-\alpha)x \cdot \alpha x + (1-\alpha)^2 xy$$

$$xy > \alpha^2 xy + \alpha y^2 - \alpha^2 y^2 + \alpha x^2 - \alpha^2 x^2 + xy - \alpha^2 xy$$

$$xy > xy + \alpha(1-\alpha)y^2 + \alpha(1-\alpha)x^2$$

we know that $\alpha(1-\alpha)y^2 + \alpha(1-\alpha)x^2 > 0$, so the above is wrong.

so $xy \leq xy + \alpha(1-\alpha)y^2 + \alpha(1-\alpha)x^2$ for $x>0, y>0, \alpha \in [0,1]$

so $f(x)$ is convex, by definition of convexity.

$$(b) f(x) = \|x\|_\infty$$

let $x, y, \alpha \in [0,1]$.

$$f(\alpha x + (1-\alpha)y) = \|\alpha x + (1-\alpha)y\|_\infty$$

by subordinate matrix infinity norm.

$$\|x\|_\infty = \max \sum_{j=1}^n |a_{ij}|$$

$$\begin{aligned} \text{so } \|\alpha x + (1-\alpha)y\|_\infty &\leq \sum_{j=1}^n (\alpha x_{ij} + (1-\alpha)y_{ij}) \\ &= \sum_{j=1}^n (\alpha x_{ij}) + \sum_{j=1}^n (1-\alpha)y_{ij}. \end{aligned}$$

$$\text{so } \|\alpha x + (1-\alpha)y\|_\infty \leq \|\alpha x\|_\infty + \|(1-\alpha)y\|_\infty = \alpha \|x\|_\infty + (1-\alpha) \|y\|_\infty$$

$$\text{so } f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

so $f(x)$ is convex, by definition of convexity.

$$(c) f(x) = x_3^3 + x_2^2 + x_1$$

$$\text{let } H = \nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_1 \partial x_3} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \frac{\partial^2 f}{\partial x_2 \partial x_3} \\ \frac{\partial^2 f}{\partial x_3 \partial x_1} & \frac{\partial^2 f}{\partial x_3 \partial x_2} & \frac{\partial^2 f}{\partial x_3^2} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 6x_3 \end{bmatrix}$$

for each part:

$$g_3(x) = x^3, \quad g_3''(x) = 6x,$$

$$g_2(x) = x^2, \quad g_2''(x) = 2,$$

$$g_1(x) = x, \quad g_1''(x) = 0.$$

$$V^T H V = [v_1 \ v_2 \ v_3] \begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 6x_3 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} = 2v_2^2 + 6v_3^2 x_3$$

for all $v, v^T \geq 0$, but for all $x, 2v_2^2 + 6v_3^2 x_3$ could be negative

so $f(x)$ is not convex, by second order condition.

$$(d) f(x) = \|Ax - b\|_2^2$$

assume A is symmetric matrix, we can set that

assume matrix A is $[m, m]$ and symmetric, assume $g(x) = Ax - b$, $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$$\begin{aligned} \text{So } H = \nabla^2 f(x) &= \left[\frac{\partial}{\partial g(x)} f(x) \frac{\partial}{\partial x} g(x) \right] \quad \text{by chain rule: } ((Ax-b)^T)^T = (2(Ax-b))' = 2A \\ &\quad \cdots \cdots \\ &\quad \left[\frac{\partial}{\partial g(x)} f(x) \frac{\partial}{\partial x} g(x) \right] \\ &= \left[2A[1,1](A[1,1] + A[2,1] + \dots + A[m,1]) \right. \\ &\quad \cdots \cdots \\ &\quad \left. - 2A[1,n](A[1,n] + A[2,n] + \dots + A[m,n]) \right] \\ &= 2A^T A \end{aligned}$$

so, for all v , we have $v^T \cdot H \cdot v = v^T \cdot 2A^T A \cdot v \geq 0$

so $H(\nabla^2 f(x))$ is positive semidefinite, so $f(x)$ is convex, by second order condition.

6. (a) Linear regression.

i. normal equations for problem (2): $\underset{\theta \in \mathbb{R}^p}{\text{minimize}} \frac{1}{2} \|X\theta - y\|_2^2$.

$$(X^T X)\theta = X^T y, \quad \theta = (X^T X)^{-1} X^T y$$

ii.

```
[13] #Question (a)
def packX(z,poly_order):
    X = np.zeros((len(z),poly_order+1))

    for i in range(poly_order+1):
        X[:, i] = np.array(z)**i

    return X

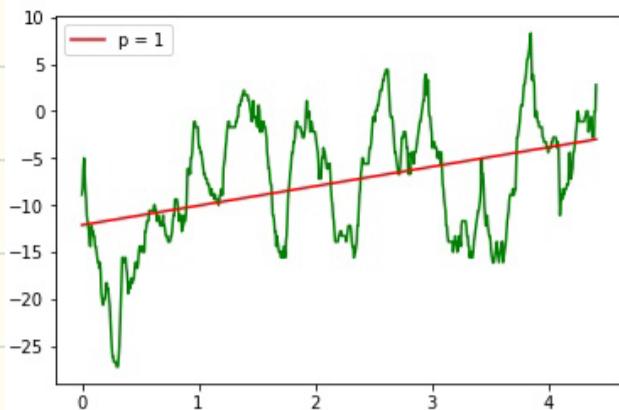
def solveLinearSystem(X,y):
    theta = np.linalg.solve(np.dot(X.T, X), np.dot(X.T, y)) #normal equations
    return theta

# TEST SCRIPT. DO NOT MODIFY!
X = packX(range(100),3)
y = np.sqrt(np.array(range(100)))
theta = solveLinearSystem(X,y)
print('Check number is', np.sum(theta))

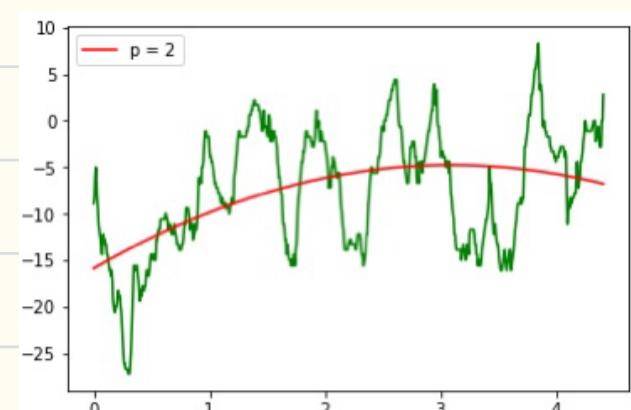
Check number is 1.341270179610586
```

Check number ≈ 1.341

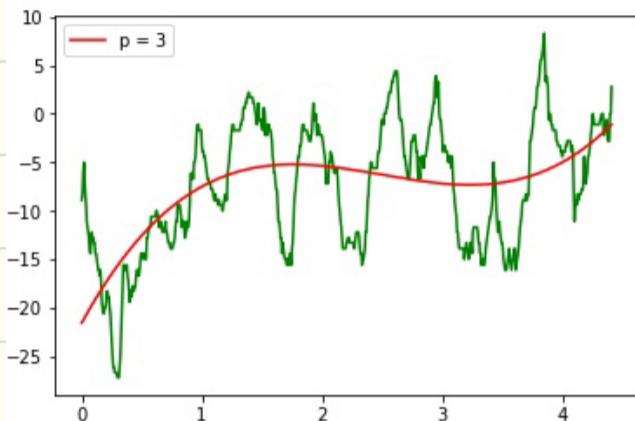
iii. $\rho = 1$.



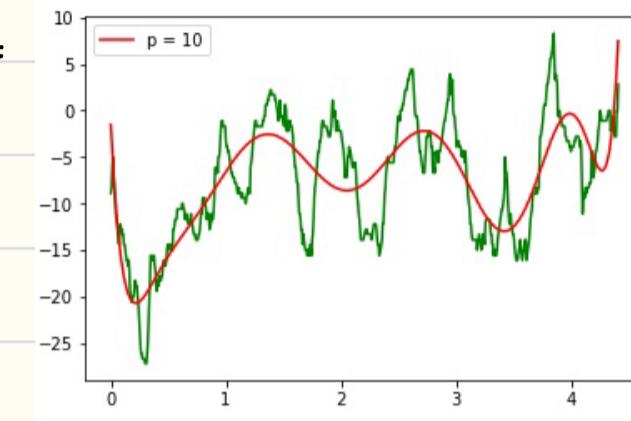
$\rho = 2$:



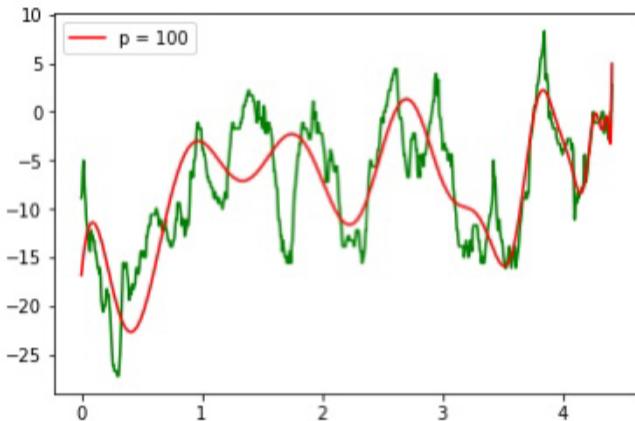
$\rho = 3$:



$\rho = 10$:



$\rho = 100$:



(b) Ridge regression.

i. normal equation for problem (3), minimize $\frac{1}{2} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \frac{\rho}{2} \|\theta\|_2^2$

$$(\mathbf{X}^T \mathbf{X} + \rho \mathbf{I}) \theta = \mathbf{X}^T \mathbf{y}$$

$$\theta = (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

ii.

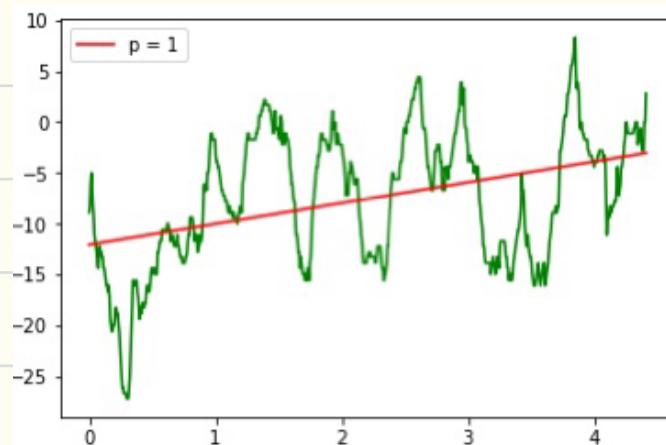
```
def solveRidgeRegressionSystem(X,y,rho):
    theta = np.linalg.solve(np.matmul(X.T, X) + rho * np.eye(3), y)
    return theta

# TEST SCRIPT. DO NOT MODIFY!
X = packX(range(100),3)
y = np.sqrt(np.array(range(100)))
theta = solveRidgeRegressionSystem(X,y,1)
print('Check number is', np.sum(theta))

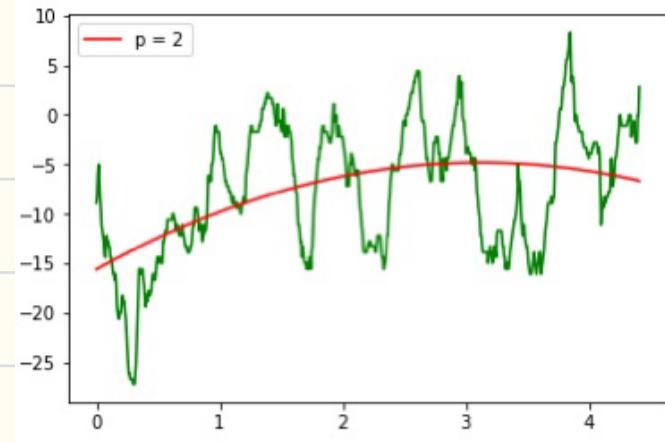
Check number is 1.2061712965226425
```

Check number is 1.206.

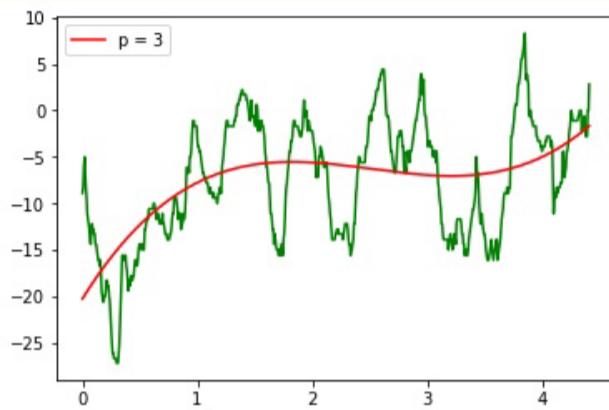
iii. $p=1$:



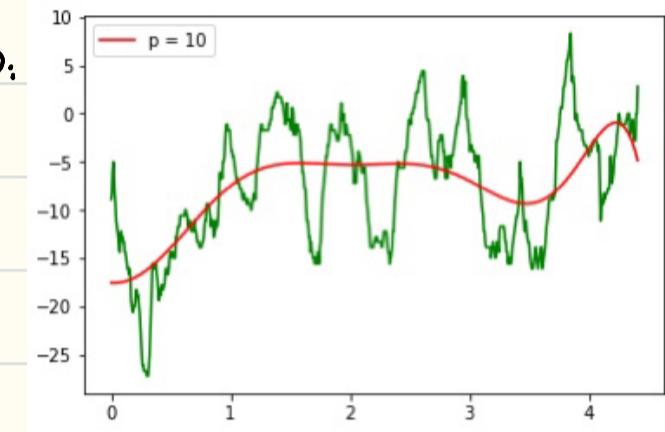
$p=2$:



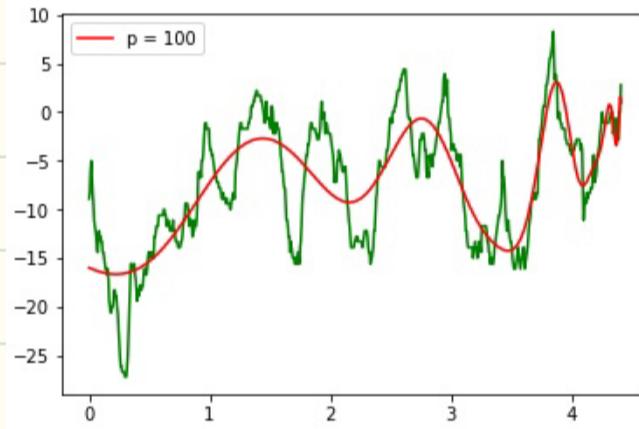
$p=3$:



$p=10$:



$p=100$:



(c) conditioning.

```
[ ] pLst1 = [1, 2, 5, 10]

▶ for p in pLst1:
    for rho in [0, N, N*10, N*100]:
        X = packX(weeks_after_start, p)
        Q = np.dot(X.T, X) + rho*np.identity(X.shape[1])
        KQ = np.linalg.cond(Q)
        print('P' + ' || ' + str(p) + ' || ' + 'rho' + ' || ' + str(rho) + ' || ' + 'K(Q)' + ' || ' + str(KQ))
```

P		1		rho		0		K(Q)		32.472521666398684
P		1		rho		742		K(Q)		6.754235401156341
P		1		rho		7420		K(Q)		1.6887590213669572
P		1		rho		74200		K(Q)		1.0702597386922308
P		2		rho		0		K(Q)		1476.3924507440943
P		2		rho		742		K(Q)		79.11371158831943
P		2		rho		7420		K(Q)		9.202205647238008
P		2		rho		74200		K(Q)		1.8243450966088195
P		5		rho		0		K(Q)		730851302.3775175
P		5		rho		742		K(Q)		271693.24960615137
P		5		rho		7420		K(Q)		27179.318088769916
P		5		rho		74200		K(Q)		2718.922773636268
P		10		rho		0		K(Q)		7.157429647449161e+18
P		10		rho		742		K(Q)		397031762841.3795
P		10		rho		7420		K(Q)		39703178289.463005
P		10		rho		74200		K(Q)		3970317849.5276527

```

y_pre = np.dot(X, theta)

mean_squared_error = np.sum((y_pre - dewtemp) ** 2)/N

print('P' + ' || ' + str(p) + ' || ' + 'rho' + ' || ' + str(rho) + ' || ' + 'mean'

```

▷ P || 1 || rho || 0 || mean square error || 36.38217361905994
P || 1 || rho || 742 || mean square error || 58.80283013948351
P || 1 || rho || 7420 || mean square error || 78.34463001407917
P || 1 || rho || 74200 || mean square error || 96.31928655068175
P || 2 || rho || 0 || mean square error || 33.53185241449123
P || 2 || rho || 742 || mean square error || 57.61883781375108
P || 2 || rho || 7420 || mean square error || 77.25615112057318
P || 2 || rho || 74200 || mean square error || 86.84731522758038
P || 5 || rho || 0 || mean square error || 27.0671045152648
P || 5 || rho || 742 || mean square error || 57.23483324805714
P || 5 || rho || 7420 || mean square error || 71.9915827133517
P || 5 || rho || 74200 || mean square error || 77.48731086603902
P || 10 || rho || 0 || mean square error || 16.707135186130255
P || 10 || rho || 742 || mean square error || 55.38371919899731
P || 10 || rho || 7420 || mean square error || 71.3923471023423
P || 10 || rho || 74200 || mean square error || 76.7444564683475

(d)

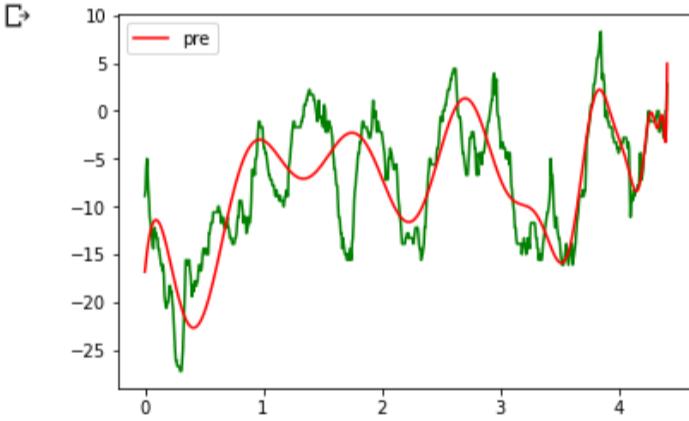
```

pM = 10
rhoM = 0
X = packX(weeks_after_start, p)
theta = solveRidgeRegressionSystem(X, dewtemp, rhoM)
y_pre = np.dot(X, theta)

plt.plot(weeks_after_start, dewtemp, c = 'green')
plt.plot(weeks_after_start, y_pre, label = 'pre', c = 'red')

plt.legend()
plt.show()

```



I pick $(p=10, \rho=0)$, now the mean square erro is the minimum (≈ 16.707).

Less erro means our predicate will be more accurate.

$$7. (a) f(\theta) = -\frac{1}{m} \sum_{i=1}^m \log(\sigma(y_i x_i^T \theta)) \quad \sigma(s) = \frac{1}{1+e^{-s}}$$

$$f_i(\theta) = \log(\sigma(y_i x_i^T \theta))$$

$$\text{Gradient: } \nabla f(\theta) = -\frac{1}{m} \sum_{i=1}^m \nabla f_i(\theta)$$

$$\nabla f_i(\theta) = (\sigma(y_i x_i^T \theta) - 1) \cdot y_i x_i$$

$$\nabla f(\theta) = -\frac{1}{m} \sum_{i=1}^m (\sigma(y_i x_i^T \theta) - 1) y_i x_i$$

(b)

```

def getLossFunction(theta):
    y = np.copy(ytrain)
    y[y== -1] = 0
    y = y.reshape((m, 1))
    h = sigmoid(X.dot(theta))
    res = ((1 / m) * X.T.dot(h - y))
    res = res.reshape(-1)
    # print(res.shape)
    return res

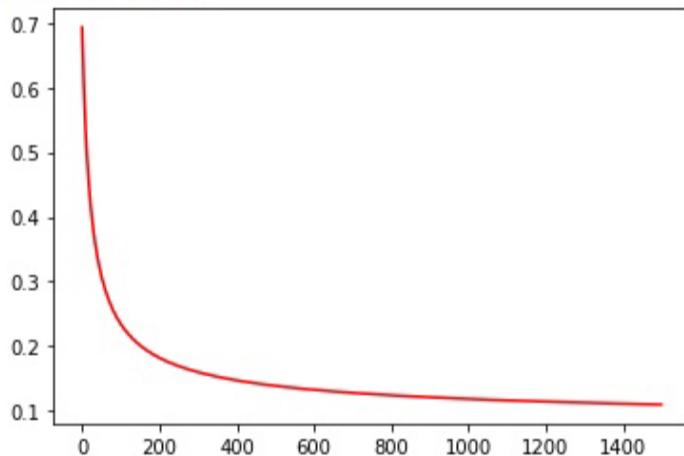
# TEST SCRIPT. DO NOT MODIFY!
theta = np.linspace(-.1,.1,n)
print('Check number is', getLossFunction(theta),np.sum(getGradient(theta)))

```

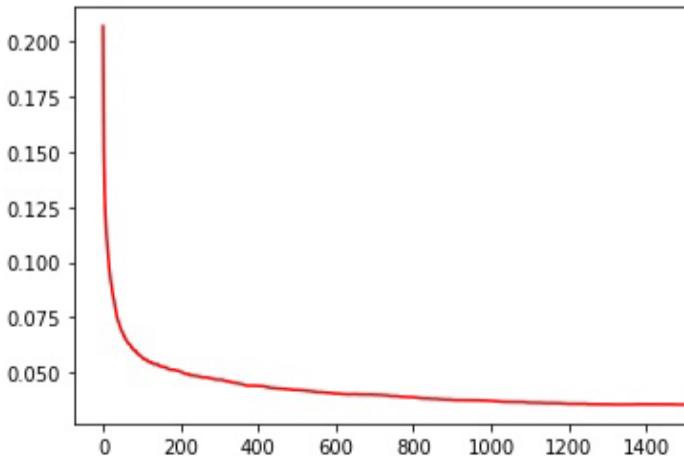
Check number is 45.19215648734918 12343.176947604472

- check number $\approx 45.192 \quad 12343.177$

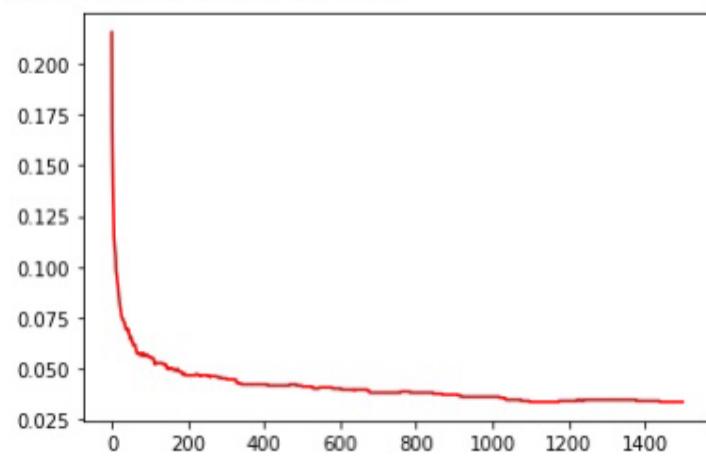
- training loss:



- training misclassification rate:



- test misclassification rate:



final train accuracy : 0.9645,
final test accuracy : 0.9663.

(c)

```
#Question(c)
def getStochGradient(theta, minibatch):
    thetaC = np.zeros(n)

    for i in minibatch:
        thetaC = thetaC + (sigmoid(ytrain[i]) * np.dot(Xtrain[i], theta))

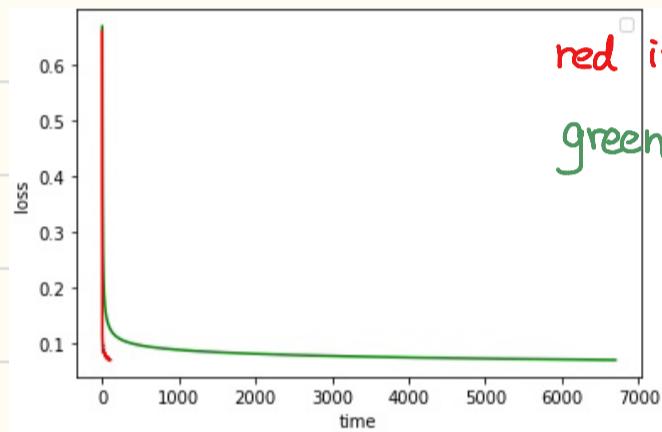
    thetaC = thetaC / len(minibatch)
    return thetaC

# TEST SCRIPT. DO NOT MODIFY!
theta = np.linspace(-.1,.1,n)
print('Check number is',np.sum(getStochGradient(theta,[1,4,1,1]),axis=0))
```

Check number is 5803.5

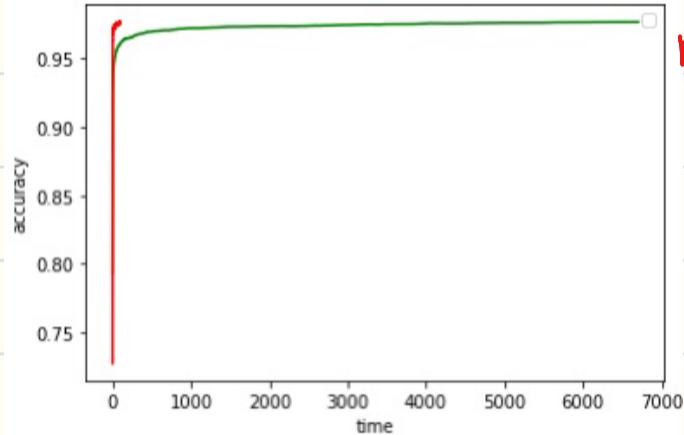
check number is 5803.5.

(d) objective loss:



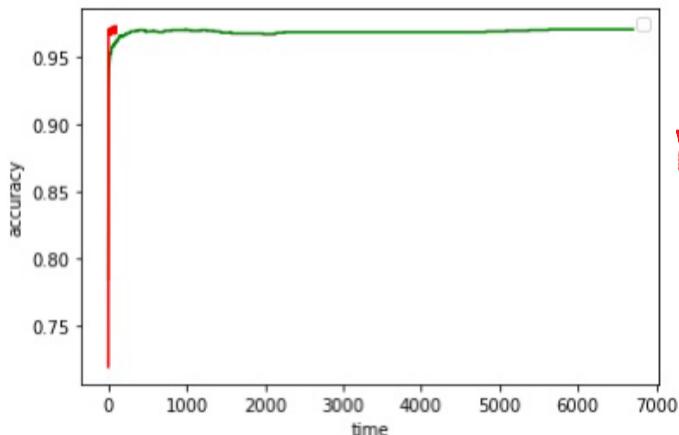
red is stochastic gradient descent
green is gradient descent

train accuracy:



red is stochastic gradient descent
green is gradient descent

test accuracy:



red is stochastic gradient descent
green is gradient descent

By running 50000 iterations, stochastic gradient descent is much slower than gradient descent. (more than 30 mins difference)

But the test accuracy of stochastic gradient descent is a bit higher than gradient descent. (about 0.002 difference)