

CS 189 HW 2

Question 1) Honor Code

1 Honor Code

1. List all collaborators. If you worked alone, then you must explicitly state so.
2. Declare and sign the following statement:

"I certify that all solutions in this document are entirely my own and that I have not looked at anyone else's solution. I have given credit to all external sources I consulted."

Signature : 

While discussions are encouraged, *everything* in your solution must be your (and only your) creation. Furthermore, all external material (i.e., *anything* outside lectures and assigned readings, including figures and pictures) should be cited properly. We wish to remind you that the consequences of academic misconduct are *particularly severe*!

Question 2 : Probability Potpourri

1) (+) Semidefinite: $n \times n$ symmetric Matrix where all of its eigenvalues are non-negative (greater or equal to 0)

$$\hookrightarrow z^T A z \geq 0 ; \text{ any vector } z$$

$$\sum_{i,j} = \text{cov}(z_i, z_j)$$

\hookrightarrow def semidefinite!

$$z^T \cdot \Sigma \cdot z ; z \in \mathbb{R}^n$$

$\underbrace{\quad}_{z \text{ is a Real number}}$



$$z^T \cdot E[(z - \mu)(z - \mu)^T] \cdot z \geq 0$$

$$\hookrightarrow E \left[\underbrace{z^T \cdot (z - \mu)}_{\text{Linear}} \underbrace{(z - \mu)^T \cdot z}_{\text{Expectation}} \right] \geq 0$$

$$\hookrightarrow E[(z - \mu) \cdot z]^2 \geq 0$$

* Since the values are being squared, the expectation is that all values will be at least 0.

$$2) X: \text{Archers hits target winds} \quad P(X) = \frac{4}{10}$$

$$Y: \text{Archers hits target hot winds} \quad P(Y) = \frac{3}{10}$$

$$Z: \text{both or none} \quad P(Z) = \frac{3}{10}$$

$$i) P(Z \cap X) = \frac{3}{10} \cdot \frac{4}{10} = \frac{12}{100}$$

$$ii) P(\text{1st shot hit}) =$$

$$\hookrightarrow = \frac{3}{10} \text{ wind} + \frac{4}{10} \text{ hit} = \frac{3}{10} \left(\frac{4}{10} \right) + \frac{7}{10} \left(\frac{3}{10} \right)$$

$$P(\text{1st shot hit}) = \frac{61}{100}$$

iii) $P(\text{Hits exactly 1 time in 2 trials})$

Binomial Distribution: k successes in n trials

$$\begin{aligned} &= \binom{2}{1} p(\text{Hit})^1 p(\text{Miss})^1 \\ &= 2 \cdot \left(\frac{3}{10} \left(\frac{4}{10} \right) + \frac{7}{10} \left(\frac{2}{10} \right) \right) \cdot \left(\frac{3}{10} \left(\frac{4}{10} \right) + \frac{7}{10} \left(\frac{2}{10} \right) \right) \\ &= 2 \cdot \frac{61}{100} \cdot \frac{31}{100} \\ &= 0.4758 \end{aligned}$$

$$\hookrightarrow P(\text{Hits exactly 1 time in 2 trials}) = 0.4758$$

$$\text{iv)} P(\text{no wins} | \text{miss}) = \frac{P(\text{miss} | \text{no wins}) P(\text{no wins})}{P(\text{miss})}$$

$$P(\text{no wins} | \text{miss}) = \frac{\frac{3}{10} \cdot \frac{7}{10}}{1 - \frac{61}{100}} = \frac{21}{100} \times \frac{100}{39}$$

$$P(\text{no wins} | \text{miss}) = \frac{7}{13}$$

$$3) E[x] = \int x \cdot f(x) dx ; x \geq 0$$

Givent: Our area Δ restricted by radii $\frac{1}{\sqrt{3}}, 1, \sqrt{2}$ f.r

Split expression level

$$E[x] = 4 \cdot P(x \leq \frac{1}{\sqrt{3}}) + 3 \cdot P(\frac{1}{\sqrt{3}} < x \leq 1) + 2 \cdot P(1 < x \leq \sqrt{2})$$

$$P(x \leq \frac{1}{\sqrt{3}}) = \int_0^{\frac{1}{\sqrt{3}}} f(x) dx ; f(x) = \frac{2}{\pi(1+x^2)}$$

$$= \frac{2}{\pi} \left[\arctan(x) \right]_0^{\frac{1}{\sqrt{3}}} ; \frac{d}{dx} \arctan = \frac{1}{1+x^2}$$

$$P(\frac{1}{\sqrt{3}} < x \leq 1) = \int_{\frac{1}{\sqrt{3}}}^1 f(x) dx$$

$$= \frac{2}{\pi} \left[\arctan(x) \right]_{\frac{1}{\sqrt{3}}}^1 ; \frac{d}{dx} \arctan = \frac{1}{1+x^2}$$

$$P(1 < x \leq \sqrt{3}) = \int_1^{\sqrt{3}} f(x) dx ; f(x) = \frac{2}{\pi(1+x^2)}$$

$$= \frac{2}{\pi} \left[\arctan(x) \right]_1^{\sqrt{3}}$$

$$E[x] = 4 \cdot P(x \leq \frac{1}{\sqrt{3}}) + 3 \cdot P(\frac{1}{\sqrt{3}} < x \leq 1) + 2 \cdot P(1 < x \leq \sqrt{3})$$

$$= \frac{2}{\pi} \left[4 \left(\arctan\left(\frac{1}{\sqrt{3}}\right) - \arctan(0) \right) + 3 \left(\arctan(1) - \arctan\left(\frac{1}{\sqrt{3}}\right) \right) + 2 \left(\arctan(\sqrt{3}) - \arctan(1) \right) \right]$$

$$= \frac{13}{6}$$

↪ $E[x] = \frac{13}{6}$

4)

$$\text{entropy}(x) = H(x) = - \sum_{x \in Z} p(x) \ln(p(x))$$

i) $X \sim \text{Bernoulli Distribution } (p)$

$$\begin{cases} p(X=1) = p \\ p(X=0) = 1-p \end{cases}$$

$\Rightarrow H(x) = - \left[p(x) \ln(p(x)) + (1-p(x)) \ln(1-p(x)) \right]$

$$H'(x) = - \left[\ln(p(x)) \cdot 1 + p(x) \cdot \frac{1}{p(x)} - \ln(1-p(x)) - \frac{1-p(x)}{1-p(x)} \right]$$

$$= - \left[\ln(p(x)) + 1 - \ln(1-p(x)) - 1 \right]$$

$$= -\ln(p(x)) + \ln(1-p(x))$$

$$H''(x) = -\frac{1}{p(x)} + \frac{-1}{1-p(x)}$$

$$= - \left(\frac{1}{p(x)} + \frac{1}{1-p(x)} \right)$$

Since $p(x) \geq 0$, This value

\Rightarrow always positive

$$\underline{H''(x) \leq 0}$$

$H(x)$ concave in p if $-H(x)$ convex function of p

ii) Sample space \geq w/ n states:

Discrete Dist

Rand Var X

→ Overall possible PMF's, entropy $H(x)$ maps to uniform Dist.

$$H(X) = - \sum_{i=1}^n p_i \ln(p_i) ; \text{ all } p \text{ equal } p_i = \frac{1}{n} \text{ for all } i$$

Log Lagrange Multiplier:

Strategy: finding local maxima/minima subject to equation constraints

$$\text{constraint } \sum_{i=1}^n p_i = 1$$

$$\hookrightarrow L = - \sum_{i=1}^n p_i \ln(p_i) + \lambda \left(\sum_{i=1}^n p_i - 1 \right)$$

$$\frac{\partial L}{\partial p_i} = -(\ln(p_i) + 1) + \lambda = 0$$

$$\hookrightarrow \ln(p_i) = \lambda - 1 \rightarrow p_i = e^{\lambda-1}$$

Plug into prob function

$$\hookrightarrow \sum_{i=1}^n e^{\lambda-1} = p_i = 1 ; \text{ identical terms}$$

$$n e^{\lambda-1} = 1$$

$$e^{\lambda-1} = p_i = \frac{1}{n} ; \text{ for all } i$$

$$\hookrightarrow H(x) = - \sum_{i=1}^n \frac{1}{n} \ln\left(\frac{1}{n}\right) = n\left(\frac{1}{n}\right) \ln\left(\frac{1}{n}\right)$$

$$= \ln(n)$$

$$\hookrightarrow H(x) = \ln(n) \text{ max entropy}$$

Question 3 : Linear Algebra Review

1) a)

3 elementary row op:

1) interchange two rows / columns

$$\begin{bmatrix} I_n & 0 \\ 0 & A \cdot B \end{bmatrix}$$

2) Multiplication of row/column by nonzero numbers

3) Multiplication of row/column by nonzero number and add result to other row/column



Left multiplies
1st row by A
and add to
2nd row

$$\begin{bmatrix} I_n & 0 \\ A & A \cdot B \end{bmatrix}$$

right multiplies 1st column
by B and subtract from
2nd column

$$\begin{bmatrix} I_n & B \\ A & 0 \end{bmatrix}$$

↓ swap column 1 with
column 2

$$\begin{bmatrix} B & I_n \\ 0 & A \end{bmatrix}$$

b) Fin) Upper/Lower Bounds Ranks ($A \cdot B$)

Prove $\text{rank}(A) + \text{rank}(B) - n \leq \text{rank}(A \cdot B)$

def Ranks Matrix: dimension of vector space generated by its columns
maximum number of linearly independent columns of A

initial matrix ranks = final matrix ranks

because elementary row operations don't impact matrix rank

$$\text{rank} \left(\begin{bmatrix} I_n & 0 \\ 0 & A \cdot B \end{bmatrix} \right) = \text{rank} \left(\begin{bmatrix} B & I_n \\ 0 & A \end{bmatrix} \right)$$

Prove $\text{rank } A + \text{rank } B - n \leq \text{rank}(AB)$

$$n + \text{rank}(AB) = \text{rank} \begin{pmatrix} I_n & 0 \\ 0 & AB \end{pmatrix} \geq \text{rank } A + \text{rank } B$$

\hookrightarrow Since AB results from $A \cdot B$, $\text{col space}(AB) \subseteq$ intersection $\text{col space } A \cap \text{col space } B$

Thus:

$$\text{rank}(AB) \leq \min(\text{rank } A, \text{rank } B)$$

c) $\text{rank}(A) = n \quad \det(M) \neq 0$

$\hookrightarrow \text{col space}(A)$ spans n linearly independent columns implying that the dimensions of col space are at least n

$\hookrightarrow \text{row space}(A)$ consists of n independent row vectors implying that its dimension is also at least n

* Since rank of matrix is number of linearly independent rows/cols, this confirms standard fact

d) Relationship between $\text{col space } A^\top$ & $\text{null space } (A^\top)$

Establish that $\text{rank}(A^\top A) = \text{rank}(A)$

$$\hookrightarrow \text{rank}(A^\top A) \leq \min(\text{rank}(A^\top), \text{rank}(A))$$

Since $\text{rank}(A^\top) = \text{rank}(A)$

$\hookrightarrow \text{rank}(A^\top A) \leq \text{rank}(A)$

Fundamental theorem linear algebra:

$$\text{Null space } (A^\top) \perp \text{column space } (A)$$

Suppose $\text{rank}(A^\top A) < \text{rank}(A)$

\hookrightarrow Implies: $A^\top A v = 0$; v nonzero vector

$\hookrightarrow A v = 0 \rightarrow$ contradicts Fundamental theorem linear algebra

Thus: $\text{rank}(A^\top A) = \text{rank}(A)$

c) $\underbrace{A\mathbb{R}^n}$ and $\underbrace{A^T A \mathbb{R}^n}$

$\hookrightarrow \mathbb{R}^n$ = input space of A

Applying A to all \mathbb{R}^n results in columns space of A

$$A\mathbb{R}^n = \text{col}(A)$$

Output of A spans entire rows of A

\hookrightarrow transformation of entire space by $A^T A$

$$A^T A \mathbb{R}^n = \text{Row space}(A)$$

$A^T A$ symmetric matrix that projects vectors onto the row space of A preserves dimensions

2) $A \in \mathbb{R}^{n \times n}$ symmetric matrix

a) $\forall x \in \mathbb{R}^n, x^T A x \geq 0$

$$v^T A v = \lambda v^T v = \lambda \|v\|^2; \lambda \text{ eigenvalue of } A \text{ corresponding eigenvector } v$$

Thus far now the conclusion

b) From part a) $\rightarrow \lambda \|v\|^2 \geq 0$

$$\text{Thus } \lambda \geq 0$$

c) Eigen decomposition of A : $A = V \Lambda V^T$; Λ diagonal matrix entries of eigenvalues of A

$\hookrightarrow V = \sqrt{\Lambda}$

$$\hookrightarrow A = V \Lambda V^T$$

$$3) Q_A(x) = x^T A x ; A = I \text{ and } x \in \mathbb{R}^n$$

Introduction knowledge : $\left\{ \begin{array}{l} A = I ; x^T A x = 1 \text{ represents circle or sphere in 2D/3D} \\ A \text{ diagonal w/ distinct val, } c_i \text{ represents an ellipse} \\ A \text{ multi 1, line or flat degenerate shape} \end{array} \right.$

a) $\underbrace{A=2I}_{\hookrightarrow} x \in \mathbb{R}^3, Q_A(x) = 32$

$$\hookrightarrow \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix} \xrightarrow{x^T A x} 2(x_1^2 + x_2^2 + x_3^2) = 32$$

* Sphere at origin with radius $\sqrt{16} = 4$

b) $A = a a^T ; a \in \mathbb{R}^3 \text{ nonzero vector}, x \in \mathbb{R}^3, Q_A(x) = 25$

$$\hookrightarrow \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{bmatrix} a_1 & a_2 & a_3 \end{bmatrix}$$

$$\hookrightarrow x^T A x = (x^T a)(a^T x) = (a^T x)^2$$

$$(a^T x)^2 = 25$$

$$a^T x = \pm 5$$

* 2 parallel planes in \mathbb{R}^3 perpendicular to a

located 15/|a|| away from the origin

c) $A = \begin{bmatrix} 5 & 1 \\ 1 & 5 \end{bmatrix} x \in \mathbb{R}^2, Q_A(x) = 24$

$$\hookrightarrow \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 5 & 1 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 5x_1 + x_2 \\ x_1 + 5x_2 \end{bmatrix} = 5x_1^2 + x_1 x_2 + x_1 x_2 + 5x_2^2$$

$$5x_1^2 + 2x_1 x_2 + 5x_2^2 = 24$$

* tilted ellipse centered at origin

d)

$$5x_1^2 + 2x_1x_2 + 5x_2^2 = 0$$

* Only solution is the origin or $x_1, x_2 = 0$

No geometric shape, simply a point on the origin

4) Frobenius inner product 2 matrices $A, B \in \mathbb{R}^{m \times n}$:

$$\langle A, B \rangle = \text{trace}(A^T B) = \sum_{i=1}^m \sum_{j=1}^n A_{i,j} B_{i,j}$$

Frobenius norm of matrix:

$$\|A\|_F = \sqrt{\langle A, A \rangle} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{i,j} B_{i,j}}$$

a) $x^T A y = \langle A, xy^T \rangle$ if $x \in \mathbb{R}^n, y \in \mathbb{R}^m, A \in \mathbb{R}^{m \times n}$

$$\hookrightarrow \sum_{i=1}^m \sum_{j=1}^n A_{i,j} \cdot x_i \cdot y_j$$

b) Using 2A:

$$A = UU^T \text{ and } b = VV^T$$

$$\begin{aligned} \text{trace}(AB) &= \text{trace}(\underbrace{UU^T}_{\text{symmetric}} VV^T) \\ &= \text{trace}(V^T U V V^T) \\ &= \text{trace}((V^T U)^2) \geq 0 \end{aligned}$$

Since $\text{trace} M \geq 0$ since it's sum of all eigenvalues
or since the trace (AB) is essentially the trace of a PSD matrix

c) PSD matrix norm $\lambda_{\max}(A)$

$$\hookrightarrow \lambda_{\max}(A) I_n - A \succeq 0$$

$$\text{prove in 4b)} \text{trace}((\lambda_{\max}(A) I_n - A) B) \geq 0$$

$$\hookrightarrow \text{trace}(AB) \leq \lambda_{\max}(A) \text{trace}(I_n B)$$

$$\hookrightarrow \text{trace}(I_n B) \leq \sqrt{n} \|B\|_F$$

Question 4b: Matrix / Vector Calculus

$$1) \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = A \quad \text{Gradient of } \sin(A_{11}^2 + e^{A_{11} + A_{22}}) + x^T A y$$

↳ To find the gradient, calculate the derivative:

$$\frac{d}{d A_{11}} \sin(A_{11}^2 + e^{A_{11} + A_{22}}) + x^T A y$$

$$\frac{d}{d A_{11}} \sin(A_{11}^2 + e^{A_{11} + A_{22}}) = \cos(A_{11}^2 + e^{A_{11} + A_{22}}) \cdot (2A_{11} + e^{A_{11} + A_{22}})$$

$$\frac{d}{d A_{22}} \sin(A_{11}^2 + e^{A_{11} + A_{22}}) = \cos(e^{A_{11} + A_{22}}) \cdot (e^{A_{11} + A_{22}})$$

↓

$$\begin{bmatrix} \cos(A_{11}^2 + e^{A_{11} + A_{22}}) \cdot (2A_{11} + e^{A_{11} + A_{22}}) & 0 \\ 0 & \cos(e^{A_{11} + A_{22}}) \cdot (e^{A_{11} + A_{22}}) \end{bmatrix}$$

$$\frac{d}{d A} x^T A y = \frac{d}{d A} \langle A, x^T y \rangle$$

$$= \frac{2}{d A} \sum_{i=1}^m \sum_{j=1}^n A_{ij} \cdot x_i^T \cdot y_j$$

$$= \frac{1}{d A} x_1 (A_{11} y_1 + A_{12} y_2) + x_2 (A_{21} y_1 + A_{22} y_2)$$

$$\frac{d}{d A_{11}} (x^T A y) = x_1 y_1, \quad \frac{d}{d A_{12}} (x^T A y) = x_1 y_2$$

$$\frac{d}{d A_{21}} (x^T A y) = x_2 y_1, \quad \frac{d}{d A_{22}} (x^T A y) = x_2 y_2$$

$$\nabla_A f(A) =$$

$$\begin{bmatrix} \cos(A_{11}^2 + e^{A_{11} + A_{22}}) \cdot (2A_{11} + e^{A_{11} + A_{22}}) + x_1 y_1 & x_1 y_2 \\ x_2 y_1 & \cos(e^{A_{11} + A_{22}}) \cdot (e^{A_{11} + A_{22}}) + x_2 y_2 \end{bmatrix}$$

$$2) \text{ a)} \alpha = \sum_{i=1}^n y_i \ln B_i ; \quad y, B \in \mathbb{R}^n$$

$$\frac{\partial \alpha}{\partial b_i} = \frac{y_i}{B_i}$$

$$\text{b)} Y = A_p + b ; \quad b \in \mathbb{R}^n \quad p \in \mathbb{R}^m \quad A \in \mathbb{R}^{m \times n}$$

$$\frac{\partial y_i}{\partial p_j} = A_{i,j}$$

$$\text{c)} \left(\frac{\partial z}{\partial x} \right)_{ii} = \frac{\partial z_i}{\partial x_i} = \sum_{j=1}^m \frac{\partial z_i}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

$$\hookrightarrow \frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

$\frac{\partial z}{\partial x}$ matrix product of other 2 Jacobians

$$1) \frac{\partial y^T z}{\partial x_i} = \sum_{j=1}^m \frac{\partial y_j}{\partial x_i} z_j + \frac{\partial z_j}{\partial x_i} y_j$$

$$\nabla_x y^T z = \frac{\partial z}{\partial x}^T z + \frac{\partial z}{\partial x}^T y$$

$$3) \text{ prove } f(x) - f(x^*) \leq \frac{p}{2}$$

$$f(x) = f(x^*) + \nabla f(x^*)^T (x - x^*) + \frac{1}{2} (x - x^*)^T H(x^*) (x - x^*)$$

Taylor series, for some x' between x^* and x

x^* global min $\rightarrow \nabla f(x^*) = 0$!

$$f(x) - f(x^*) = \frac{1}{2} (x - x^*)^T H(x^*) (x - x^*)$$

$H(x^*)$ positive semidefinite so max eigenvalue is 1:

$$\underbrace{\lambda_{\max}(H(x^*))}_{2} \leq 1$$

Quod erat demonstrandum:

$$(x - x^*)^T H(x^*) (x - x^*) \leq \lambda_{\max}(H(x^*)) \|x - x^*\|^2$$

$$\hookrightarrow f(x) - f(x^*) \leq \frac{1}{2} \|x - x^*\|^2$$

x belongs in spherical region

$$X = \{x \mid \|x - x^*\|^2 \leq n\}$$

$$\hookrightarrow f(x) - f(x^*) \leq \frac{D}{2}$$

thus:

$$f(x) - f(x^*) \leq \frac{D}{2} \quad \forall x \in X$$

Question 6: Properties of Normal Dist

1) Prove: $E[e^{\lambda x}] = e^{\sigma^2 \lambda^2 / 2}$; Hint MGF

$$\begin{aligned} \hookrightarrow E[e^{\lambda x}] &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{\lambda x} e^{-x^2/2\sigma^2} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{\lambda x - x^2/2\sigma^2} e^{-z^2/2} dz \\ &= e^{\sigma^2 \lambda^2 / 2} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(z - \lambda\sigma)^2/2} dz \\ &= e^{\sigma^2 \lambda^2 / 2} \end{aligned}$$

2)

Gaussian RV are independent if not correlated:

$$E[u_x v_x] = E\left(\left[\sum_{i=1}^n u_i x_i\right] \left[\sum_{i=1}^n v_i x_i\right]\right)$$

$$= \underbrace{\sum_{i=1}^n u_i v_i}_{\langle u, v \rangle} E[x_i^2]$$

$$\langle u, v \rangle = 0$$

$$= 0$$

$\Rightarrow u_x$ and v_x are independent. However, if u, v are not iid, this can come probably with the correlation value. We are only able to derive $\text{cov} = 0$ due to iid given. Thus, w/o this fact, we cannot say that u_x and v_x are independent.

Question 7: Multivariate Normal Distribution

1) $E[x] = \mu$

$$\hookrightarrow f(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

$$E[x] = \int_{\mathbb{R}^d} x \cdot f(x) dx$$

Change of variables

def 2: $\Sigma^{-1/2}(x - \mu)$; standardizes x so $z \sim N(0, I_d)$
 $x = \mu + \Sigma^{1/2} z$

$$\hookrightarrow f_z(z) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2} z^T z}$$

$$\hookrightarrow E[x] = \int_{\mathbb{R}^d} (\mu + \Sigma^{1/2} z) f_z(z) dz$$

$$= \mu + \Sigma^{1/2} \int_{\mathbb{R}^d} z f_z(z) dz$$

$$\underbrace{\quad}_{=0 \text{ w/ symmetry of standard multivariate}}$$

$$\text{normal dist. } E[z] = 0$$

$$E[x] = \mu$$

$$2) \text{Var}(x) = \Sigma$$

$$\text{Var}(x) = \text{Cov}(x, x) = E[(x - \mu)(x - \mu)^T] = E[xx^T] - \mu\mu^T$$

$$x - \mu = \sum_{i=1}^n z_i ; \text{ change of variable: } x - \mu = \sum_{i=1}^n z_i$$

$$\text{Var}(x) = E[(\sum_{i=1}^n z_i)(\sum_{i=1}^n z_i)^T]$$

$$= E[(\sum_{i=1}^n z_i z_i^T \sum_{i=1}^n z_i)] ; z \sim N(0, I_2)$$

$$= \sum_{i=1}^n E[z z^T] \sum_{i=1}^n$$

$$= \sum_{i=1}^n I_2 \sum_{i=1}^n$$

$$= \Sigma$$

$$\therefore \text{Var}(x) = \Sigma$$

Question 8: Gradient Descent

optimization problem: $\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x - b^T x ; A \in \mathbb{R}^{n \times n} \text{ PSD w/ } \lambda_{\min}(A) < \lambda_{\max}(A)$

1) Objective is convex \rightarrow optimum stationary point of objective:

$$Ax - b = 0 ; A \text{ invertible}$$

$$\hookrightarrow \text{optimizer: } x^* = A^{-1}b$$

$$2) x^{k+1} = x^k - (Ax^k - b)$$

$$3) x^k - x^* = (I - A)(x^{k-1} - x^*)$$

$$x^k - x^* = x^{k-1} - (Ax^{k-1} - b) - x^*$$

$$= (I - A)x^{k-1} + b - x^*$$

$$= (I - A)x^{k-1} - (I - A)x^*$$

$$= (I - A)(x^{k-1} - x^*)$$

$$4) \|Ax\|_2 \leq \lambda_{\max}(A) \|x\|_2$$

$$\hookrightarrow \|Ax\|_2 = x^T A^T x \leq \lambda_{\max}(A)^2$$

If $x \neq 0$, then:

$$\|A(x/\|x\|_2)\|_2 \leq (\lambda_{\max}(A))^2$$

which proves $\|Ax\|_2 \leq \lambda_{\max}(A) \|x\|_2$ after some manipulation

$$5) \|x^{(k)} - x^*\|_2 \leq p \|x^{(k-1)} - x^*\|_2$$

Since $\lambda_{\max}(A) < 1$; $I - A > 0$:

$$\|(I - A)(x^{(k-1)} - x^*)\|_2 \leq \lambda_{\max}(I - A) \|x^{(k-1)} - x^*\|_2$$

$$p = \lambda_{\max}(I - A) = 1 - \lambda_{\min}(A)$$

$$\|x^{(k)} - x^*\|_2 = \|(I - A)(x^{(k-1)} - x^*)\|_2 \leq p \|x^{(k-1)} - x^*\|_2$$

b) recursion formula:

$$\|x^{(k)} - x^*\|_2 \leq p^k \|x^0 - x^*\|_2$$

$$\text{Conditioning on } p^k \|x^0 - x^*\|_2 \leq \varepsilon$$

$$\hookrightarrow \log(p^k) \cdot \log(\|x^0 - x^*\|_2) \leq \log(\varepsilon)$$

$$\hookrightarrow \log(p) \cdot \log(\|x^0 - x^*\|_2) \leq \log(\varepsilon)$$

$$\text{Is } \log(p) \leq \log\left(\frac{\|x^0 - x^*\|_2}{\varepsilon}\right)$$

$$\text{Is } \leq \frac{1}{\log(p)} \log\left(\frac{\|x^0 - x^*\|_2}{\varepsilon}\right)$$

Question 5: Linear Neural Networks

Applies multivariate chain rule to simple Neural Net \rightarrow Linear Neural Network

- Only univariate Linear regression / decision func
- $n \times d$ design matrix X
 - ↳ Each row of X trains point
 - ↳ X n training points w/ d features
- $n \times k$ matrix Y
 - ↳ Each row of Y set of labels for each X training point

Goal: Learn $k \times d$ matrix W or weight w^T such that:

$$Y \approx Xw^T$$

Find W minimizes loss func:

$$RSS(W) = \|Xw^T - Y\|_F^2$$

Instead of optimizing W over $ls \times d$ space, write W in terms

as product of multiple matrices

Linear Neural Net:

$$W = \nu(w_L, w_{L-1}, \dots, w_1) = w_L w_{L-1} \dots w_1$$

↳ ↳
 Matrix multiplication row vectors $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{L-1}$, matrix
 map

Weight vector: $\theta = (w_L \dots w_1) \in \mathbb{R}^d$

$$\nu(\theta) = \nu(w_L \dots w_1)$$

↳
 column vector

Goal: Find θ that min RSS

$$J(\theta) = RSS(\nu(\theta))$$

Using Gradient Descent

$$1) \quad g = \nabla_w RSS(w) \quad k \times 2 \text{ matrix}$$

$$\hookrightarrow RSS(w) = \sum_{i=1}^n \sum_{j=1}^k ((xw^T)_{ij} - y_{ij})^2 ; \|A\|_F^2 = \sum_{i,j} A_{ij}^2$$

$$(xw)^T_{ij} = \sum_{m=1}^2 x_{im} w_{jm}$$

$$RSS(w) = \sum_{i=1}^n \sum_{j=1}^k \left(\sum_{m=1}^2 x_{im} w_{jm} - y_{ij} \right)^2$$

$$g_{pq} = \frac{\partial RSS(w)}{\partial w_{pq}}$$

$$\hookrightarrow \frac{\partial}{\partial w_{pq}} RSS(w) = 2 \sum_{i=1}^n \sum_{j=1}^k \left(\sum_{m=1}^2 x_{im} w_{jm} - y_{ij} \right) \cdot \frac{\partial}{\partial w_{pq}} \sum_{m=1}^2 x_{im} w_{jm} ; \quad j=p \quad m=q \text{ for } w_{jm}$$

$$\hookrightarrow \frac{\partial}{\partial w_{pq}} \sum_{m=1}^2 x_{im} w_{jm} = x_{ip} \delta_{pq}$$

$$\hookrightarrow g_{pq} = 2 \sum_{i=1}^n (xw^T - y)_{ip} x_{ip} \quad \text{Summation form}$$

$$R = xw^T - y$$

$$G = 2R^T X$$

$$\hookrightarrow G = 2(xw^T - y)^T X$$

$$G = 2(xw^T x - y^T x) \quad \text{Gradient form}$$

2) Directional Derivative of RSS(w)

$$RSS'_{\Delta w}(w) = \lim_{\epsilon \rightarrow 0} \frac{RSS(w + \epsilon \Delta w) - RSS(w)}{\epsilon}$$

first order approximation

$$RSS(w + \epsilon \Delta w) \approx RSS(w) + \epsilon \nabla_w RSS(w) \cdot \Delta w + O(\epsilon^2)$$

$$RSS'_{\Delta w}(w) = \nabla_w RSS(w) \cdot \Delta w ; \quad \nabla_w RSS(w) = 2(xw^T - y)^T X$$

$$\hookrightarrow RSS'_{\Delta w}(w) = 2(xw^T - y)^T X \cdot \Delta w$$

$$\langle A, B \rangle_F = \sum_{i,j} A_{ij} B_{ij}$$

$$= \text{Tr}(A^\top B)$$

$$R_{\Delta w}^F(w) = 2 \text{Tr}((xw^\top - y)^\top x \Delta w)$$

$$= 2 \text{Tr}((xw^\top - y)^\top x \Delta w)$$

$$\langle A, B \rangle_F = \text{Tr}(A^\top B)$$

$$\hookrightarrow R_{\Delta w}^F(w) = 2 \langle xw^\top - y, x \Delta w^\top \rangle_F$$

3)

$$w = \nu(w_L, w_{L-1}, \dots, w_1)$$

$$= w_L \dots w_1$$

$$N'_{\Delta \theta}(\theta) ; \theta = (w_L \dots w_1) \quad \Delta \theta = (\Delta w_L, \Delta w_{L-1}, \dots, \Delta w_1)$$

$$N'_{\Delta \theta}(\theta) = \sum_{j=1}^L w_j^> \Delta w_j w_j^<; \quad w_j^> = w_L w_{L-1} \dots w_{j+1}$$

$$w_j^< = w_{j-1} w_{j-2} \dots w_1$$

$$w_L^> = I_{d_L} \quad (\text{identity matrix of size } d_L \times d_L)$$

$$w_1^< = I_{d_0} \quad (\text{identity matrix of size } d_0 \times d_0)$$

$$R_{\Delta w}^F(w) = \lim_{\epsilon \rightarrow 0} \frac{R(w + \epsilon \Delta w) - R(w)}{\epsilon} ; \quad N(\theta) = w_L w_{L-1} \dots w_1$$

$$N'_{\Delta \theta}(\theta) = \sum_{j=1}^L w_L \dots w_{j+1} \Delta w_j w_{j-1} \dots w_1$$

$$w_j^> = w_L \dots w_{j+1} \quad w_j^< = w_{j-1} \dots w_1$$

$$\hookrightarrow N'_{\Delta \theta}(\theta) = \sum_{j=1}^L w_j^> \Delta w_j w_j^<$$

41 Multivariate Chain Rule $\mathcal{J}'_{\theta\theta}(\theta)$

$$\hookrightarrow \mathcal{J}'_{\theta\theta}(\theta) = \mathbb{E}_{w'_{\Delta w}(w)} \cdot \mathcal{N}'_{\theta\theta}(\theta)$$

Multivariate Chain Rule

$$\frac{\partial}{\partial x} f(g(x)) = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\hookrightarrow Q2: \mathbb{E}_{w'_{\Delta w}(w)} = \mathbb{E} \langle xw^T - y, x\Delta w^T \rangle_F$$

$$\hookrightarrow Q3: \mathcal{N}'_{\theta\theta}(\theta) = \sum_{j=1}^L w_j^T \Delta w_j w_j^T$$

↳ Chain Rule: $\mathcal{J}'_{\theta\theta}(\theta) = \mathbb{E}_{w'_{\theta\theta\theta\theta}(w)} \mathcal{N}'_{\theta\theta}(\theta)$

$$\mathcal{J}'_{\theta\theta}(\theta) = \mathbb{E} \langle xw^T - y, x \left(\sum_{j=1}^L w_j^T \Delta w_j w_j^T \right)^T \rangle_F$$

Linearity of inner product:

$$\mathcal{J}'_{\theta\theta}(\theta) = \mathbb{E} \sum_{j=1}^L \langle xw^T - y, x (w_j^T)^T \Delta w_j^T (w_j^T)^T \rangle_F$$

$$\mathcal{J}'_{\theta\theta}(\theta) = \mathbb{E} \sum_{j=1}^L \langle w_j^T (xw^T - y)^T x w_j^T, \Delta w_j \rangle_F$$

$$\hookrightarrow \nabla_{\theta} \mathcal{J}(\theta) = (2w_1^T (xw^T - y)^T x w_L^T, \dots, 2w_s^T (xw^T - y)^T x w_s^T, \dots, 2w_L^T (xw^T - y)^T x)$$

S) $\nabla_{\theta} S(\theta)$

$$\nabla_{\theta} S(\theta) = (2w_1^T (xw^T - y)^T xw_1^L, \dots, 2w_s^T (xw^T - y)^T xw_s^L, \dots, 2w_t^T (xw^T - y)^T x)$$

Q2: $RSS'_{\theta w}(w) = 2 \langle xw^T - y, x\Delta w^T \rangle_F$

Q3: $\nabla_{\theta\theta} S(\theta) = \sum_{j=1}^L w_j^T \Delta w_j w_j^L$

Q4:

$$\nabla_{\theta} S(\theta) = (2w_1^T (xw^T - y)^T xw_1^L, \dots, 2w_s^T (xw^T - y)^T xw_s^L, \dots, 2w_t^T (xw^T - y)^T x)$$

$$\nabla_{w_i} S(\theta) = 2w_i^T (xw^T - y)^T xw_i^L ; \quad w^T = (w_1^T \ w_2^T \ \dots \ w_L^T)$$

$$\nabla_{w_i} S(\theta) = 2w_i^T (xw^T - y)^T xw_i^L$$