

HW02

Who I had discussions with:

Went to physical OTs and got help from various students.

Also got help/helped Madeline (Don't know her last name) Also looked in the CSISq discord for help usually when stuck. Some thing for the EDStem posts on the homework.

"I certify that all solutions are entirely my own words and that I have not looked at another student's solutions. I have given credit to all external sources I consulted."

X Erin

$$1) \quad 1) \quad 1\{a \geq b\} = \begin{cases} 1 & \text{if } a \geq b \\ 0 & \text{if } a < b \end{cases}$$

← Indicator function

$$P(A) = E[1\{A\}]$$

$$A = X \geq t$$

$$P(X \geq t) = E[1\{X \geq t\}]$$

$$1\{X \geq t\} \leq \frac{X}{t}$$

$\left\{ \begin{array}{l} 1\{a \geq b\} \leq \frac{a}{b} \\ \text{because if } a \geq b \text{ then} \\ \text{RHS} = 1 \text{ which is the} \\ \text{highest value the indicator} \\ \text{funct. can output.} \end{array} \right.$

$$P(X \geq t) \leq E\left[\frac{X}{t}\right] \quad \text{Note: } E\left[\frac{X}{t}\right] = \frac{1}{t} E[X]$$

because $\frac{1}{t}$ is a constant

$$\boxed{P(X \geq t) \leq \frac{E[X]}{t}}$$

2)

$$X = (1\hat{\mu} - \mu)^2$$

$$(P(|1\hat{\mu} - \mu| \geq t))^2 \leq \frac{E[(1\hat{\mu} - \mu)^2]}{(t)^2}$$

$$(P(|1\hat{\mu} - \mu| \geq t))^2 \leq \frac{1}{n(t)^2}$$

Taken from previous part
 $P(X \geq t) \leq E[X]/t$
 but $x = 1\hat{\mu} - \mu$

$$\boxed{P(|1\hat{\mu} - \mu| \geq t) \leq \frac{1}{\sqrt{n}t}}$$

$$2) \quad 1) \quad x^T \underline{\Sigma} x \geq 0 \quad \longleftarrow \text{A definition of PSD} \quad (1)$$

$$x^T E[(2-\mu)(2-\mu)^T] x \geq 0$$

$$E[\underline{x^T (2-\mu) (2-\mu)^T x}] \geq 0$$

The two underlined products are identical we will call w

$$E[w] \geq 0$$

$$E[w^2] \geq 0$$

Something squared is always greater than 0,
so covariance matrix
satisfies a PSD equation.

$$2) \quad i) 0.3 \cdot 0.4 = 0.12 \quad \boxed{12\%}$$

$$ii) 0.4 \cdot 0.7 = 0.28 = \boxed{28\%}$$

$$iii) 0.28 \cdot 0.7 = 0.2016 = \boxed{20\%}$$

$$iv) 0.6 = \boxed{60\%}$$

$$3) \quad \int_0^{1/\sqrt{3}} \frac{4}{\pi(1+x^2)} dx + \int_{1/\sqrt{3}}^1 \frac{2}{\pi(1+x^2)} dx + \int_1^{\sqrt{3}} \frac{2}{\pi(1+x^2)} dx$$

$$\frac{8}{\pi} [\arctan(x)]_0^{1/\sqrt{3}} + \frac{6}{\pi} [\arctan(x)]_{1/\sqrt{3}}^1 + \frac{4}{\pi} [\arctan(x)]_1^{\sqrt{3}}$$

$$\frac{4}{3} + \frac{1}{2} + \frac{1}{3} = \boxed{\frac{13}{6}}$$

3) is a)

$$\begin{bmatrix} I_n & 0 \\ 0 & AB \end{bmatrix}$$

$$\begin{bmatrix} 0 & I_n \\ AB & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 \\ B & 1 \end{bmatrix}$$

$$\begin{bmatrix} B & I_n \\ 0 & AB \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ 0 & AB \end{bmatrix}$$

$$\begin{bmatrix} B & I_n \\ 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 0 & A \end{bmatrix}$$

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

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

b)

Upper Bound:

rank of $\begin{bmatrix} B & I_n \\ 0 & A \end{bmatrix}$ must be the same as the rank of

$\begin{bmatrix} I_n & 0 \\ 0 & AB \end{bmatrix}$ which this rank is n (Because $\text{rank } I_n = n$)
 $+ \text{rank}(AB)$ and rank of first matrix is $\text{rank}(A) + \text{rank}(B)$,
 we have: $n + \text{rank}(AB) = \text{rank}(A) + \text{rank}(B)$

$$\boxed{\text{rank}(AB) \geq \text{rank}(A) + \text{rank}(B) - n}$$

Lower Bound:

Since the final matrix dimension is $m \times p$, this means the row space of AB is in the subspace of A and the column the subspace of B . This means the rank is dependent on A and B , so that means it cannot exceed the lower one of the two, leaving us with $\text{rank}(AB) \leq \min(\text{rank } A, \text{rank } B)$

c) $\det(M) \neq 0$ Therefore M must be full rank with a row space dimension of r and column space dimension of r .

d) Nullspace of $A^T A$ is defined as $A^T A x = 0$. If the rank $(A^T A)$ is less than $\text{rank}(A)$, then that means there are less lin. ind. columns in $A^T A$ than A . Because of rank nullity theorem this means if rank is smaller then nullspace must be bigger which implies there must be a non-trivial solution in there. If there is a non-trivial solution to $A^T A x = 0$ then that means it must be in the column space of A and nullspace of A^T which are orthogonal, which means they cannot be shared together. This causes a contradiction.

e) $A R^n$ is the matrix sized to fit the R^n space. Specifically for the row space since R^n is the $n \times 1$ matrix.

$A^T A R^n$ is the same as mentioned before but instead it applies to the row and column space

2) $A \rightarrow B$:

$$Ax = \lambda x$$

$$x^T Ax \geq 0$$

$$x^T \lambda x \geq 0 \quad \leftarrow \lambda \text{ is just a number}$$

$$\lambda x^T x \geq 0 \quad \leftarrow x^T x \geq 0 \text{ due to squared nature}$$

So the only way the above inequality is true is if $\lambda \geq 0$

$B \rightarrow C$:

$$A = \lambda D \lambda^T$$

- λ contains eigenvectors
- D diagonal matrix with eigenvalues
- Diagonalized the matrix

$$D = D^{1/2} D^{1/2} \quad \text{Aka Singular values}$$

$$A = \lambda D^{1/2} D^{1/2} \lambda^T$$

$$A = \lambda D^{1/2} (\lambda D^{1/2})^T$$

$$U = \lambda D^{1/2}$$

$$\boxed{A = UU^T}$$

$C \rightarrow A:$

$$x^T A x$$

$$x^T (U U^T) x$$

$$x^T U U^T x$$

$$x^T U (U^T U)^T$$

$$\omega = x^T U$$

$$\omega \omega^T \geq 0$$

Therefore inequality holds true.

3) a)

$$\langle A, x y^T \rangle$$

$$\text{trace}(A^T x y^T) \quad \text{order does not matter}$$

$$\text{trace}(x A^T y)$$

$$\text{trace}(x^T A y)$$

Note on dimensions: $|x| \times |A| \times |y| = |x|$

$$|x| \times |A| \times |y| = |x|$$

So $x^T A y$ is a scalar

The trace of a scalar is itself.

b) For PSD we know $x^T A x \geq 0$
 $x^T B x \geq 0$

therefore since A, B are both are positive, then the trace as a result would also be positive.

c) Create a diagonal matrix with the diagonal entries are

$$\langle A, B \rangle = \text{trace}(A B)$$

$x^T A B x \geq 0$ Scaling by a constant $\lambda_{\max}(A)$ we will be able to ensure holds true

$$\underbrace{\|A\|_F \|B\|_F}_{\langle A, B \rangle} \leq \sqrt{n} \lambda_{\max}(A) \|B\|_F$$

$$4) r(AA^T, U) = \frac{U^T (AA^T) U}{U^T U}$$

$$\frac{U^T (A^T A) U}{U^T U} = r(A^T A, U)$$

$$\frac{U^T (AA^T) U}{U^T U} = \frac{U^T (A^T A) U}{U^T U}$$

As a result of this, the singular value of A as shown in SVD decomposition $U \Sigma V^T$ the Σ would be represented as $\max_{U \in \mathbb{R}^{m \times n}, \|U\|_F = 1} U^T A V$

$$4) \quad 1) \quad \begin{bmatrix} 2A_{11} \cos(A_{11}^2 + e^{A_{11} + A_{22}}) e^{A_{11} + A_{22}} + x^T A y & \cos(A_{11}^2) e^{A_{11} + A_{22}} + x^T A y \\ \cos(A_{11}^2) e^{A_{11} + A_{22}} + x^T A y & \cos(A_{11}^2 + e^{A_{11} + A_{22}}) e^{A_{11} + A_{22}} + x^T A y \end{bmatrix}$$

$$2) \quad a) \quad \sigma_{\max}(A) = \sqrt{\lambda_{\max}(A^T A)}$$

$$b) \quad \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$$

$$3) \quad a) \quad \sum_{i=1}^n \frac{\partial \alpha}{\partial B_i} y_i \ln B_i$$

$$\sum_{i=1}^n \frac{y_i}{B_i}$$

$$b) \quad \frac{\partial y_i}{\partial p_j} = A$$

$$c) \quad k \times m \quad m \times n \quad k \times n$$

$$\begin{bmatrix} \frac{d_2}{d\gamma} & \frac{d\gamma}{dx} \end{bmatrix}$$

$$d) \quad \nabla_{\gamma} \mathcal{L} + \gamma^T \nabla_{\gamma} \mathcal{L}$$

$$4) a) \quad \sum_{i=1}^m \sum_{j=1}^n (X W^T - Y)^2 \quad \leftarrow \text{with summation}$$

$$(X W^T - Y)^2 \quad \leftarrow \text{without summation}$$

$$b) \quad \langle \text{RSS}(w), \Delta W \rangle$$

$$c) \quad \mu'_{\Delta\theta}(\theta) = \Delta\theta \cdot \mu(\theta)$$

$$= \sum_{j=1}^L w_{\rightarrow j} \Delta w_j w_{\leftarrow j}$$

$\mu(\theta)$ term
is irrelevant due
to summation

d) The chain rule for scalar functions is as follows:

$$\frac{d}{dx} f(g(x))|_{x=x_0} = \frac{d}{dx} f(y)|_{y=g(x_0)} \cdot \frac{d}{dx} g(x)|_{x=x_0}$$

applying this onto RSS and μ , we get

$$\boxed{\frac{d}{d\theta} \|XW^T - Y\|_F^2 \cdot \frac{d}{d\theta} \mu(\theta)}$$

e) $J(\theta) = \text{RSS}(\mu(\theta))$

$$= \|XW^T - Y\|_F^2$$

$$\sum_{j=1}^L \langle XW^T - Y \rangle^2 \quad (1)$$

$$\mu(\theta) = \sum_{j=1}^L W_{\rightarrow j} \Delta W_j W_{\leftarrow j} \quad (2)$$

Putting these all together 1, 2, 3 gives us.

$$W = \mu(W_L, W_{L-1}, \dots, W_2, \dots) \quad (3)$$

$$\boxed{\begin{aligned} & (2(\mu(\theta))X^T - Y^T)XW_{\leftarrow L}^T \dots \\ & 2W_{\rightarrow j}^T (\mu(\theta)X^T - Y^T)X \\ & W_{\leftarrow j}^T \dots 2W_{\rightarrow j}^T (\mu(\theta)X^T - Y^T) \\ & X) \end{aligned}}$$

$$5) 1) \quad E[e^{\lambda X}] = e^{\sigma^2 \lambda^2 / 2}$$

$$MGF = \int_{-\infty}^{\infty} e^{\lambda x} f(x) dx$$

$$M_X(\lambda) = E[e^{\lambda X}]$$

PDF = Since Normal Gaussian, then

$$\frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{x^2}{2\sigma^2}}$$

$$\text{Which means } MGF = \int_{-\infty}^{\infty} e^{\lambda x} \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{x^2}{2\sigma^2}} dx = e^{\sigma^2 \lambda^2 / 2}$$

2)

$$P(X \geq t) \leq \frac{E[X]}{t}$$

$$X = e^{\lambda X} \quad t = e^{\lambda t}$$

$$\frac{E[e^{\lambda X}]}{e^{\lambda t}}$$

$$P(X \geq t) \leq \frac{e^{\frac{\sigma^2 \lambda^2}{2}}}{e^{\lambda t}}$$

$$\lambda = \frac{t}{\sigma^2}$$

$$P(X \geq t) \leq \frac{\exp\left(\frac{\sigma^2 \left(\frac{t^2}{\sigma^4}\right)}{2}\right)}{\exp\left(\frac{t}{\sigma^2} + \right)}$$

$$\frac{\exp\left(\frac{t^2}{2\sigma^2}\right)}{\exp\left(\frac{t}{\sigma^2}\right)}$$

$$\frac{t^2}{2\sigma^2} - \frac{2t}{2\sigma^2} = \frac{-t^2}{2\sigma^2}$$

$$\boxed{e^{-\frac{t^2}{2\sigma^2}}}$$

3) As $n \rightarrow \infty$ the concentration inequality will go closer and closer to the average

4) Yes u_x and v_x are both independent because of X being i.i.d. However, if X becomes independent, but not identically distributed, then u_x and v_x become dependant as now the values in X are dependant on some i , when before, they were not, it was just $N(0,1)$.

6) 1)

$$\int_{\mathbb{R}^d} x f(x; \mu, \Sigma) dx$$

$$z = \Sigma^{-1/2} (x - \mu)$$

$$\Sigma^{1/2} z = x - \mu$$

$$\Sigma^{1/2} z + \mu = x$$

$$\int_{\mathbb{R}^d} (\Sigma^{1/2} z + \mu) f(\Sigma^{1/2} z + \mu; \mu, \Sigma) dz$$

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

$$\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot \exp\left[-\frac{1}{2} (\Sigma^{1/2} z)^T \Sigma^{-1} (\Sigma^{1/2} z)\right]$$

$$11 \cdot \exp\left[-\frac{1}{2} (z^T \underbrace{\Sigma^{1/2} \Sigma^{-1} \Sigma^{1/2}}_I) z\right]$$

$$11 \cdot \exp\left[-\frac{1}{2} (z^T z)\right]$$

$$\int_{\mathbb{R}^d} \left(\underbrace{\Sigma^{1/2} z}_{\Sigma^{1/2} \cdot 0} + \underbrace{\mu}_{\mu \cdot 1} \right) dz$$

$\underbrace{\qquad\qquad\qquad}_{\pi}$

$$6) 2) Z = \Sigma^{1/2} (X - \mu)$$

$$\text{Var}(Z) = E[ZZ^T] - E[Z]E[Z]^T$$

For $E[ZZ^T]$:

$$\begin{aligned} E[ZZ^T] &= E[\Sigma^{1/2} (X - \mu) (\Sigma^{1/2} (X - \mu))^T] \\ &= \Sigma^{1/2} E[(X - \mu)(X - \mu)^T] \Sigma^{1/2} \\ &= \Sigma^{1/2} E[(X - \mu)(X - \mu)] \Sigma^{1/2} \\ &\quad \underbrace{\Sigma^{1/2} \Sigma \Sigma^{1/2}}_{= \Sigma} = \Sigma \end{aligned}$$

For $E[Z]E[Z]^T$:

$$E[Z]E[Z]^T$$

$$E[\Sigma^{1/2} (X - \mu)] \cdot E[(X - \mu) \Sigma^{1/2}]$$

$$\Sigma^{1/2} E[X - \mu] \cdot E[X - \mu] \Sigma^{1/2}$$

$$\Sigma^{1/2} (E[X] - \mu) \cdot (E[X] - \mu) \cdot \Sigma^{1/2}$$

$$\Sigma^{1/2} (\mu - \mu) \cdot (\mu - \mu) \cdot \Sigma^{1/2} = 0$$

Therefore $\sum_1 - 0 = \boxed{\sum_1 = \text{Var}(x)}$

$$7) 1) \text{Gradient} = Ax - b$$

Local min @ Gradient = 0

$$Ax - b = 0$$

$$+b \quad +b$$

$$\cancel{A^{-1}A}x = \vec{A}b$$

$$\boxed{x^* = A^{-1}b} \leftarrow \text{Opt. mizer}$$

2) Next pt. = curr. pt. - Step · gradient

$$\boxed{x^{k+1} = x^k - 1 \cdot (Ax^k - b)}$$

$$3) x^k - (Ax^k - b) = x^{k+1}$$

$$x^k - Ax^k + b$$

$$x^k (I - A) + b$$

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

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

$$x^{k+1} - x^* = x^k - (Ax^k + Ax^*) - (\cancel{A^{-1}A}x^*)$$

$$x^k - x^* - (Ax^k + x^*)$$

$$x^{(k)} - Ax^{(k)} - x^* - x^*$$

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

This is expressed as
 $k+1$ and k , but
 you could just
 do -1 on all
 the k exponents
 to get the formula
 in the question

$$Av = \lambda v$$

$$4) \|Ax\|^2 = (Ax)^T Ax = x^T A^T Ax$$

Rayleigh quotient

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

$$= x^T A^2 x$$

$$= x^T A^2 x$$

?
~~only~~
 if symmetric
 Yes PSD
 means
 symmetric

$$\sqrt{x^T A^2 x} \leq \sqrt{\lambda_{\max}(A^2) \cdot \|x\|_2^2}$$

$$\boxed{\|Ax\|^2 \leq \lambda_{\max}(A) \cdot \|x\|_2^2}$$

$$5) \quad x^k - x^* = (I - A)(x^{k-1} - x^*)$$

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

6) k should be

$$\boxed{\frac{\|x^k - x^*\|_2}{\rho \epsilon}}$$

References:

- <https://math.stackexchange.com/questions/326784/prove-that-the-rank-of-a-block-diagonal-matrix-equals-the-sum-of-the-ranks-of-th>
- <https://math.stackexchange.com/questions/349738/prove-operatorname-rank-a-operatorname-rank-a-for-any-a-in-m-m-times-n>
- <https://math.stackexchange.com/questions/456025/what-is-the-relationship-between-the-null-space-and-the-column-space>
- [https://en.wikipedia.org/wiki/Rank_\(linear_algebra\)](https://en.wikipedia.org/wiki/Rank_(linear_algebra))
- <https://math.stackexchange.com/questions/29072/how-is-the-column-space-of-a-matrix-a-orthogonal-to-its-nullspace>
- https://en.wikipedia.org/wiki/Frobenius_inner_product
- <https://math.stackexchange.com/questions/3516928/why-does-xtx-x2>
- https://en.wikipedia.org/wiki/Diagonalizable_matrix
- https://ocw.mit.edu/courses/res-18-011-algebra-i-student-notes-fall-2021/mit18_701f21_lect12.pdf
- <https://www.khanacademy.org/math/linear-algebra/alternate-bases/orthonormal-basis/v/lin-alg-orthogonal-matrices-preserve-angles-and-lengths>
- <https://math.stackexchange.com/questions/2958603/if-a-pdpt-does-p-have-to-be-orthogonal>
- [https://math.libretexts.org/Bookshelves/Linear_Algebra/Fundamentals_of_Matrix_Algebra_\(Hartman\)/03%3A_Operations_on_Matrices/3.01%3A_The_Matrix_Transpose](https://math.libretexts.org/Bookshelves/Linear_Algebra/Fundamentals_of_Matrix_Algebra_(Hartman)/03%3A_Operations_on_Matrices/3.01%3A_The_Matrix_Transpose)
- [https://en.wikipedia.org/wiki/Trace_\(linear_algebra\)](https://en.wikipedia.org/wiki/Trace_(linear_algebra))
- https://en.wikipedia.org/wiki/Definite_matrix
- https://en.wikipedia.org/wiki/Cauchy%E2%80%93Schwarz_inequality
- https://en.wikipedia.org/wiki/Singular_value_decomposition
- <https://math.stackexchange.com/questions/213653/show-that-norm-of-matrix-a-is-given-by-the-square-root-of-the-largest-eigenval>
- https://en.wikipedia.org/wiki/Singular_value
- https://en.wikipedia.org/wiki/Markov%27s_inequality

- https://www.stat.berkeley.edu/~mjlwain/stat210b/Chap2_TailBounds_Jan22_2015.pdf
- https://en.wikipedia.org/wiki/Concentration_inequality
- <https://en.wikipedia.org/wiki/Estimator>
- https://en.wikipedia.org/wiki/Jensen%27s_inequality
- <https://math.stackexchange.com/questions/3148049/reverse-of-jensens-inequality>
- <https://www.geeksforgeeks.org/covariance-matrix/>
- <https://en.wikipedia.org/wiki/Covariance>
- <https://math.stackexchange.com/questions/1524424/prove-that-the-average-of-iid-gaussian-random-variables-is-gaussian>
- [https://en.wikipedia.org/wiki/Rank_\(linear_algebra\)](https://en.wikipedia.org/wiki/Rank_(linear_algebra))
- <https://byjus.com/jee/rank-of-a-matrix-and-special-matrices/>
- <https://www2.math.upenn.edu/~moose/240S2013/slides7-22.pdf>
- <https://math.stackexchange.com/questions/64350/what-is-the-difference-between-the-row-space-and-the-column-space-in-linear-algebr>
- https://en.wikipedia.org/wiki/Row_and_column_spaces
- <https://stattrek.com/matrix-algebra/elementary-operations>
- <https://www.khanacademy.org/math/linear-algebra/vectors-and-spaces/null-column-space/v/introduction-to-the-null-space-of-a-matrix>
- https://en.wikipedia.org/wiki/Rank%E2%80%93nullity_theorem
- <https://www.algebrapracticeproblems.com/how-to-diagonalize-a-matrix-diagonalizable-diagonalization/>
- https://en.wikipedia.org/wiki/Moment-generating_function#:~:text=1%20For%20a%20discrete%20probability%20mass%20function%2CM%20X,%28t%29%3Dint%20%20%7B-infty%20%7D%5E%20%7Binfy%20%7De%5E...%20More%20
- https://ocw.mit.edu/courses/18-600-probability-and-random-variables-fall-2019/d68c0ee3be77495da2aeaf921f7fc081_MIT18_600F19_lec26.pdf
- <https://www.scribbr.com/methodology/independent-and-dependent-variables/>

- <https://brilliant.org/wiki/expected-value/>
- https://en.wikipedia.org/wiki/Matrix_norm
- <https://math.stackexchange.com/questions/289989/first-and-second-derivative-of-a-summation>
- https://en.wikipedia.org/wiki/Partial_derivative
- https://en.wikipedia.org/wiki/Jacobian_matrix_and_determinant
- <https://en.wikipedia.org/wiki/Gradient>
- https://en.wikipedia.org/wiki/Artificial_neural_network#:~:text=The%20simplest%20kind%20of%20feedforward%20neural%20network%20%28FNN%29,and%20the%20inputs%20is%20calculated%20at%20each%20node.
- <https://www.mathworks.com/help/deeplearning/ug/linear-neural-networks.html>
- https://en.wikipedia.org/wiki/Design_matrix
- <https://builtin.com/machine-learning/cost-function>
- https://en.wikipedia.org/wiki/Linear_least_squares
- https://en.wikipedia.org/wiki/Matrix_norm#Frobenius_norm
- <https://math.stackexchange.com/questions/1898839/relation-between-frobenius-norm-and-trace>
- [https://en.wikipedia.org/wiki/Projection_\(mathematics\)](https://en.wikipedia.org/wiki/Projection_(mathematics))
- https://math.iit.edu/~fass/477577_Chapter_10.pdf
- [https://math.libretexts.org/Bookshelves/Calculus/Calculus_\(OpenStax\)/14%3ADifferentiation_of_Functions_of_Several_Variables/14.05%3A_The_Chain_Rule_for_Multivariable_Functions](https://math.libretexts.org/Bookshelves/Calculus/Calculus_(OpenStax)/14%3ADifferentiation_of_Functions_of_Several_Variables/14.05%3A_The_Chain_Rule_for_Multivariable_Functions)