

TP 1

YUNLONG JIAO

21/01/2019

PART I: Basics of Linear Algebra and Matrix Calculus

Def. (Matrix): Let $a_{ij} \in \mathbb{R}$ *only real matrices unless o.w.s.* $i=1, \dots, m, j=1, \dots, n.$

$${}_m A_n = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

Rmk:

① ${}_n A_n$ square matrix

② ${}_n A_1 := a \in \mathbb{R}^n$ vector

Conv.

(bold print) A, B, X, U, V matrices

(bold print) $\vec{a}, \vec{b}, \vec{x}, \vec{y}$ vectors

(bold print) $\lambda, \alpha, a_i, x_i, a_{ij}$ scalars (reals)

(bold print) $O_n := [0]$ null matrix

(bold print) $I_n := \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}$ identity matrix

Def. (Matrix Multiplication)

$${}_m C_p = {}_m A_n \cdot {}_n B_p \Leftrightarrow c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$

compatible for mult.

$$i \begin{pmatrix} \cdot \\ \cdot \\ \cdot \end{pmatrix} = i \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \begin{pmatrix} | \\ | \\ | \end{pmatrix}$$

Rmk:

① ${}_m z = {}_m A_n \cdot {}_n x \Leftrightarrow z_i = \sum_{j=1}^n a_{ij} x_j$

② $\alpha = {}_1 y^T A_n \cdot {}_n x \Leftrightarrow \alpha = \sum_{j=1}^n \sum_{k=1}^m a_{kj} x_j y_k$

Prop. (Multiply partitioned matrices)

$$A = \begin{array}{c} m \\ n \end{array} \left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline A_{21} & A_{22} \end{array} \right] \begin{array}{c} n_1 \quad n_2 \\ m_1 \quad m_2 \end{array} \quad B = \begin{array}{c} n \\ p \end{array} \left[\begin{array}{c|c} B_{11} & B_{12} \\ \hline B_{21} & B_{22} \end{array} \right] \begin{array}{c} n_1 \quad n_2 \\ p_1 \quad p_2 \end{array}$$

$$AB = \begin{array}{c} m \\ p \end{array} \left[\begin{array}{c|c} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ \hline A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{array} \right] \begin{array}{c} n_1 \quad n_2 \\ p_1 \quad p_2 \end{array}$$

Def. (Basic matrix operations)

Transpose ${}_n A_m^T = [a_{ji}]$

For square matrix $\left\langle \text{Inverse } {}_n A_n^{-1} \text{ s.t. } A^{-1}A = AA^{-1} = I \right.$

Trace $\text{tr}(A) = \sum_{i=1}^n a_{ii}$

Prop. (Invertible matrices)

${}_n A_n$ invertible (A^{-1} exists)

$\Leftrightarrow Ax=0$ has only the trivial solution $x=0$.

i.e. $\ker(A) = 0$.

i.e. Columns of A are linearly independent.

Ex. ① $(AB)^T = B^T A^T$

$(AB)^{-1} = B^{-1} A^{-1}$ *always assumed invertible if $^{-1}$ exists.*

$(A^T)^{-1} = (A^{-1})^T$ *always assumed compatible for mult.*

② $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$

$\text{tr}(\lambda A) = \lambda \text{tr}(A)$

$\text{tr}(AB) = \text{tr}(BA)$

$\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$

③ $\text{tr}(AA^T) = \sum_{i,j} a_{ij}^2 = \|A\|_F^2$ Frobenius norm of matrices

Ex: (Matrix Inversion Lemma)

$$\textcircled{1} (A + \underset{\substack{n \quad m \quad m \quad n}}{UCV})^{-1} = A^{-1} - A^{-1}U(\underset{\substack{m \quad n \quad n \quad m}}{C^{-1} + VA^{-1}U})^{-1}VA^{-1}$$

$$\textcircled{2} (I + \vec{u}\vec{v}^T)^{-1} = I - \frac{uv^T}{1 + v^T u}$$

Ex: (Schur complement, Invert partitioned matrix)

Suppose $M = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$, A, D, M invertible.

Denote Schur complements by

$$M/D := A - BD^{-1}C$$

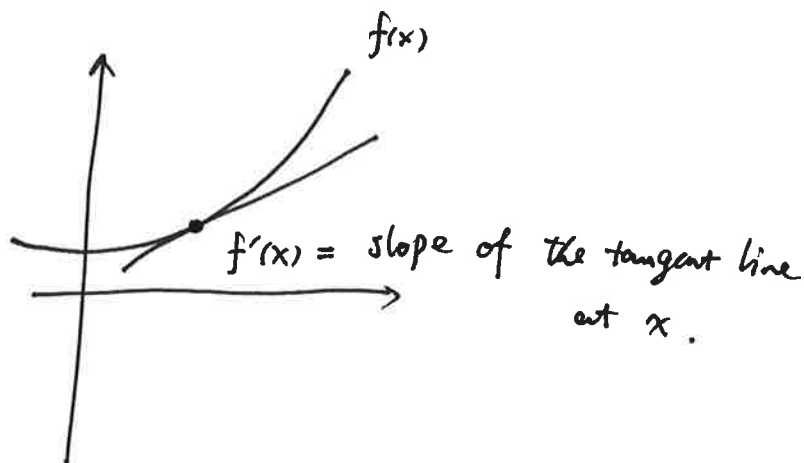
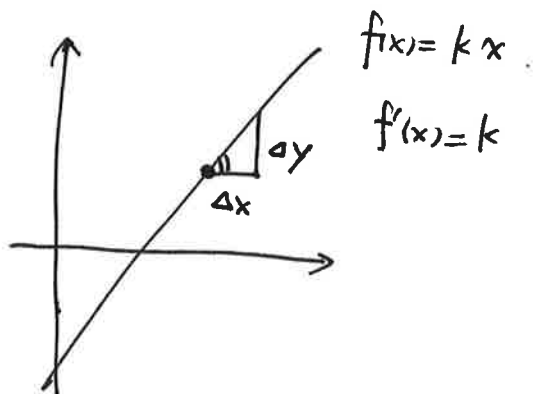
$$M/A := D - CA^{-1}B$$

Show that

$$M^{-1} = \begin{bmatrix} (M/D)^{-1} & - (M/D)^{-1}BD^{-1} \\ -D^{-1}C(M/D)^{-1} & (M/A)^{-1} \end{bmatrix}$$

Interlude: Derivative (of a function of a real variable)

$$y = f(x) \quad \frac{dy}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} (= f'(x))$$



Def. (Matrix Derivatives, Jacobian, Hessian)

① $\vec{y} = f(\vec{x})$ vector-valued function of vectors

Jacobian of \vec{f} (matrix) $\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$

② $\alpha = f(\vec{x})$ real-valued function of vectors

Jacobian of f (vector) $\frac{\partial \alpha}{\partial \vec{x}} = \begin{pmatrix} \frac{\partial \alpha}{\partial x_1} \\ \vdots \\ \frac{\partial \alpha}{\partial x_n} \end{pmatrix} := \underline{J(\alpha)}$

③ \star Hessian of f (matrix) $\frac{\partial^2 \alpha}{\partial x_i \partial x_j} = \begin{bmatrix} \frac{\partial^2 \alpha}{\partial x_1^2} & \dots & \frac{\partial^2 \alpha}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 \alpha}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 \alpha}{\partial x_n^2} \end{bmatrix} := \underline{H(\alpha)}$

③ $A = f(\alpha)$ matrix-valued function of reals

Jacobian of f (matrix) $\frac{\partial A}{\partial \alpha} = \left[\frac{\partial A_{ij}}{\partial \alpha} \right]$

④ $\alpha = f(A)$ real-valued function of matrices

Jacobian of f $\frac{\partial \alpha}{\partial A} = \left[\frac{\partial \alpha}{\partial A_{ij}} \right]$

Rmk: ① All defined by element-by-element derivatives of real functions.
 ② The order of Jacobian may depend on A , or A^T by convention.

Ex:

$$\textcircled{1} \quad \vec{y} = A\vec{x} \quad \frac{\partial \vec{y}}{\partial \vec{x}} = A$$

proof $y_i = \sum_{j=1}^n a_{ij} x_j$, $\frac{\partial y_i}{\partial x_j} = a_{ij}$

$$\textcircled{2} \quad \alpha = \vec{y}^T A \vec{x}$$

$$\frac{\partial \alpha}{\partial \vec{x}} = A^T \vec{y} \quad , \quad \frac{\partial \alpha}{\partial \vec{y}} = A \vec{x}$$

$$\textcircled{3} \quad \alpha = \vec{x}^T A \vec{x}$$

$$J_x(\alpha) = (A + A^T) \vec{x}$$

$$H_x(\alpha) = A + A^T$$

$$\textcircled{4} \quad \varepsilon = \|\vec{y} - A\vec{x}\|_2^2 = \sum_i (y_i - \sum_j a_{ij} x_j)^2$$

$$\frac{\partial \varepsilon}{\partial \vec{x}} = -2 A^T (\vec{y} - A\vec{x})$$

~~Ex~~ Ex (Chain rule)

Not a function of \vec{z} , unless specified!

$$\textcircled{1} \quad \vec{y} = A\vec{x} \quad , \quad \vec{x} \text{ is a function of } \vec{z}$$

$$\frac{\partial \vec{y}}{\partial \vec{z}} = A \frac{\partial \vec{x}}{\partial \vec{z}}$$

proof LHS: by def. $\frac{\partial \vec{y}}{\partial \vec{z}} = \left[\frac{\partial y_i}{\partial z_j} \right]$, $y_i = \sum_{k=1}^n a_{ik} x_k$

$$\Rightarrow \frac{\partial y_i}{\partial z_j} = \sum_{k=1}^n a_{ik} \frac{\partial x_k}{\partial z_j}$$

RHS: by def. $\frac{\partial \vec{x}}{\partial \vec{z}} = \left[\frac{\partial x_k}{\partial z_j} \right]$ by chain rule of diff. of real functions

by matrix mult. $\left[A \frac{\partial \vec{x}}{\partial \vec{z}} \right]_{ij} = \sum_{k=1}^n a_{ik} \frac{\partial x_k}{\partial z_j}$

$$\textcircled{2} \quad \alpha = \vec{y}^T A \vec{x} \quad , \quad \vec{x}, \vec{y} \text{ are both functions of } \vec{z}$$

$$\frac{\partial \alpha}{\partial \vec{z}} = \vec{x}^T A^T \frac{\partial \vec{y}}{\partial \vec{z}} + \vec{y}^T A \frac{\partial \vec{x}}{\partial \vec{z}}$$

Ex: ① Suppose A is a (matrix-valued) function of α

$$\frac{\partial(A^{-1})}{\partial \alpha} = -A^{-1} \frac{\partial A}{\partial \alpha} A^{-1}$$

② $\frac{\partial}{\partial A} \text{tr}(AB) = B^T$

$m \quad n \quad p$

$$\frac{\partial}{\partial A} \text{tr}(A) = I$$

compare $\left\{ \begin{array}{l} \frac{\partial}{\partial A} \text{tr}(ABA^T) = A(B+B^T) \\ \frac{\partial}{\partial B} \text{tr}(ABA^T) = A^T A \end{array} \right.$

③ $\mathcal{E} = \|X - WH\|_F^2 = \sum_i \sum_k (x_{ik} - \sum_{j \in \mathcal{J}} w_{ij} h_{jk})^2$

$$\frac{\partial \mathcal{E}}{\partial W} = -2XH^T + 2WHH^T$$

P A U S E "BREAK"

Def. (Eigenvector Equation)

$${}_n A_n \cdot \underbrace{A \frac{u_i}{\uparrow}}_{\text{eigenvector}} = \underbrace{\lambda_i}_{\text{eigenvalue}} \frac{u_i}{\uparrow} \quad i = (1, \dots, n)$$

Rmk. rank(A) = number of non-zero eigenvalues.

Def. (Symmetric Matrix)

$${}_n A_n \text{ symmetric iff } A = A^T$$

Def. (Orthogonal Matrix)

$${}_n U_n \text{ orthogonal iff } U^T = U^{-1} \left(\begin{array}{l} \Rightarrow U U^T = U^T U = I \\ \Rightarrow u_i^T u_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} = \delta_{ij} \end{array} \right) \quad | \quad 6$$

why?

Thm. (Eigendecomposition)

Any (real) symmetric matrix $A_{n \times n}$ can be decomposed:

$$A_{n \times n} = U_{n \times n} \Lambda_{n \times n} U_{n \times n}^T$$

where U is an orthogonal matrix,
whose columns are eigenvectors of A .

$\Lambda := \text{diag}(\lambda_1, \dots, \lambda_n)$ is a diagonal matrix,
whose diagonal are eigenvalues of A .

Cor.: ① Symmetric matrix A invertible $\Leftrightarrow \lambda_i \neq 0 \quad \forall i$

② $A^{-1} = U \Lambda^{-1} U^T$, $\Lambda^{-1} = \text{diag}(\lambda_1^{-1}, \dots, \lambda_n^{-1})$

③ why? $\text{tr}(A) = \sum_i \lambda_i$, $\text{tr}(A^{-1}) = \sum_i \lambda_i^{-1}$

Def.: (Positive (semi-) definite matrix)

$A_{n \times n} \succeq 0$ p.s.d. iff $x^T A x \geq 0, \forall x \in \mathbb{R}^n$.

$A_{n \times n} \succ 0$ p.d. iff $x^T A x > 0, \forall x \in \mathbb{R}^n$.

Prop. (Eigenvalues of Symmetric p.(s)d. matrices)

① $A_{n \times n} \succeq 0$ p.s.d. and Symm. \Leftrightarrow e.v. $\lambda_i \geq 0, i=1, \dots, n$

② $A_{n \times n} \succ 0$ p.d. and symm. \Leftrightarrow e.v. $\lambda_i > 0, i=1, \dots, n$ \rightarrow full-rank matrix

Thm. (Singular Value Decomposition, SVD)

Any (real) matrix A can be decomposed as

$$A = U \Sigma V^T$$

$m \quad n \quad m \quad m \quad n \quad n$

Where U, V are orthogonal matrices

whose columns are left- and right- singular vectors of A
 Σ is a "diagonal" matrix $m \times n$ with $\begin{pmatrix} \diagdown & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \diagdown \end{pmatrix}$ OR $\begin{pmatrix} \diagdown & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \diagdown \end{pmatrix}$ non-negative real numbers on diagonal, called Singular values of A , denoted by $\sigma_i \geq 0, i=1, \dots, \min(m, n)$.

Ex: For any general X show that

① $X^T X$ and $X X^T$ are symmetric, p.s.d.

② The non-zero eigenvalues of $X^T X$ and $X X^T$ are the same, that are $\{\sigma_i^2 \mid \sigma_i \neq 0, i=1, \dots, \min(m, n)\}$ where σ_i 's are singular values of X .

8

PART II: Intro to Constrained Convex Optimization

Linear Regression:

$$\mathcal{D} = \left\{ (\vec{x}_i, y_i) \right\}_{i=1}^m, \quad \vec{x}_i \in \mathbb{R}^n, y_i \in \mathbb{R}$$

$$y_i = \underbrace{\vec{w}^T \vec{x}_i}_{\text{model fit}} + \underbrace{\varepsilon_i}_{\text{error}}, \quad \vec{w} = ?$$

Squared loss: $\|\varepsilon\|^2 = \sum_{i=1}^m \varepsilon_i^2 = \sum_{i=1}^m (y_i - \vec{w}^T \vec{x}_i)^2$

$$=: \|\vec{y} - X\vec{w}\|_2^2$$

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_m^T \end{pmatrix} \in \mathbb{R}^{m \times n} \quad \vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

Empirical Risk Minimization (ERM):

$$\min_{\vec{w} \in \mathbb{R}^n} \|\vec{y} - X\vec{w}\|_2^2 := \ell(\vec{w})$$

$$\frac{\partial \ell}{\partial \vec{w}} = -2 X^T (\vec{y} - X\vec{w}) \stackrel{\text{Why? Fermat's!}}{=} 0$$

$$(X^T X) \vec{w} = X^T \vec{y}$$

If $(X^T X)$ invertible, $\hat{\vec{w}}^{\text{OLS}} = (X^T X)^{-1} X^T \vec{y}$

Rmk:

① (Ordinary) least-squares solution to linear regression
 $\hat{\vec{w}}^{\text{OLS}}$. Best Linear Unbiased Estimate (Gauss-Markov)

② $(X^T X)$ may not be invertible \rightarrow pseudo-inverse ...
 or invertible but ill-conditioned. \rightarrow condition number $\kappa = \frac{\lambda_{\max}(X^T X)}{\lambda_{\min}(X^T X)} \approx 0$
 $\hat{\vec{w}}^{\text{OLS}}$ is unbiased but can have very large variance $\propto \frac{1}{\lambda_{\min}(X^T X)}$

to alleviate the problem!

Regularization

"Ridge"

$$\min_w \|y - Xw\|_2^2$$

$$\text{s.t. } \|w\|_2^2 \leq t$$



Why?

Guaranteed by strong duality!

$$\exists \lambda > 0$$

$$\min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

dual optimal

Lagrangian

$l(w)$

$$\frac{\partial l}{\partial w} = -2X^T(y - Xw) + 2\lambda w = 0$$

$$\hat{w}^{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

Rank: $(X^T X + \lambda I)$ is always invertible.

and $\lambda_{\min}(X^T X + \lambda I) \geq \lambda$

Rank: Constrained opt \rightarrow unconstrained opt \rightarrow Fermat's

In general, what to do ???

Move to slide presentation / BREAK

(in a vector space)

Def. (Convex set) X convex set iff

$$\forall x_1, x_2 \in X, \forall t \in [0, 1], t x_1 + (1-t) x_2 \in X$$

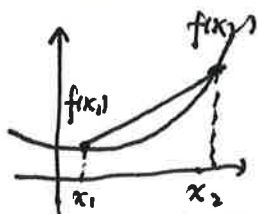
Def. (Convex function) $f: X \rightarrow \mathbb{R}$, X strictly convex set iff

$$\forall x_1, x_2 \in X, \forall t \in [0, 1], f(t x_1 + (1-t) x_2) < t f(x_1) + (1-t) f(x_2)$$

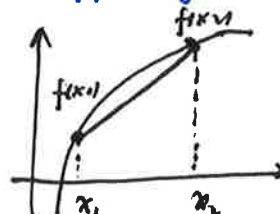
Def. (Concave function) g is concave iff $-g$ is convex.



eg: \mathbb{R}^n
ball $B_0 = \{x : \|x\|_2 \leq t\}$



eg: Convex (e^x)



Concave ($\log x$)