

Lecture 2: Numerical Analysis

Instructor: Yifan Chen

Scribes: Yifan Xu

Proof reader: Yifan Chen, Xiong Peng

Note: LaTeX template courtesy of UC Berkeley EECS dept.

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

2.1 Matrix derivation

Matrix derivation refers to the process of computing the derivative of a matrix-valued function with respect to the input matrix, or the derivative of a scalar function to the input matrix. In this section, we study the latter with the matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$, and the scalar function $f(\mathbf{X}) \in \mathbb{R}$. The derivative $\frac{\partial f}{\partial \mathbf{X}}$ can be defined using element-wise derivation:

$$\frac{\partial f}{\partial \mathbf{X}} = \left[\frac{\partial f}{\partial \mathbf{X}_{ij}} \right]_{m \times n}. \quad (2.1)$$

Computing element-wise derivatives for large matrices is complex and lacks elegance. We aim to interpret $f(\mathbf{X})$ as a function of the matrix \mathbf{X} rather than as a function of its individual elements \mathbf{X}_{ij} . Therefore, to establish a matrix-centric perspective on matrix derivation, we first revisit the derivatives we have previously learned.

In univariate calculus, the derivative maps a scalar to a scalar, and the differential is expressed as:

$$df = f'(x)dx,$$

where df is the differential, $f'(x)$ is the derivative.

For multivariable calculus, the derivative maps a scalar to a vector, and the differential is defined as:

$$df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f^T}{\partial \mathbf{x}} d\mathbf{x},$$

where the total differential df is the inner product of the gradient $\frac{\partial f}{\partial \mathbf{x}}$ and the differential vector $d\mathbf{x}$. We generalize the total differential concept from multivariate calculus to matrix calculus as follows:

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial \mathbf{X}_{i,j}} d\mathbf{X}_{i,j} = \text{Tr} \left(\frac{\partial f^T}{\partial \mathbf{X}} d\mathbf{X} \right) = \left\langle \frac{\partial f}{\partial \mathbf{X}}, d\mathbf{X} \right\rangle, \quad (2.2)$$

where $\text{Tr}(\cdot)$ represents matrix trace, which is the sum of the diagonal elements of a square matrix.

Note: Trace $\text{Tr}(\cdot)$ satisfies the property that for matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$, $\text{Tr}(\mathbf{A}^T \mathbf{B}) = \sum_{i,j} \mathbf{A}_{ij} \mathbf{B}_{ij}$, i.e., $\text{Tr}(\mathbf{A}^T \mathbf{B})$ is the inner product of matrices \mathbf{A} and \mathbf{B} .

Now we can use differential to compute derivative, we first build rules for basic differential operations.

2.1.1 Differential formulas

1. $d(\mathbf{X} + \mathbf{Y}) = d\mathbf{X} + d\mathbf{Y}$ (Addition)
 2. $d(\mathbf{XY}) = d\mathbf{X} \cdot \mathbf{Y} + \mathbf{X} \cdot d\mathbf{Y}$ (Multiplication)
 3. $d\mathbf{X}^{-1} = -\mathbf{X}^{-1}d\mathbf{XX}^{-1}$ (Inverse)
- This formula can be proven using $d\mathbf{XX}^{-1} = d\mathbf{I}$.

4. $d(\mathbf{X} \odot \mathbf{Y}) = d\mathbf{X} \odot \mathbf{Y} + \mathbf{X} \odot d\mathbf{Y}$ (Element-wise multiplication),
where \odot represents element-wise multiplication of matrices \mathbf{X} and \mathbf{Y} of the same size.
5. $d\sigma(\mathbf{X}) = \sigma'(\mathbf{X}) \odot d\mathbf{X}$, $\sigma(\mathbf{X}) = [\sigma(\mathbf{X}_{ij})]$ (Element-wise function),
where $\sigma(\mathbf{X}) = [\sigma(\mathbf{X}_{ij})]$ represents element-wise function, $\sigma'(\mathbf{X}) = [\sigma'(\mathbf{X}_{ij})]$ represents element-wise derivative.

E.g., for matrix $\mathbf{X} = \begin{bmatrix} \mathbf{X}_{11} & \mathbf{X}_{12} \\ \mathbf{X}_{21} & \mathbf{X}_{22} \end{bmatrix}$,

$$d \sin(\mathbf{X}) = d \begin{bmatrix} \sin \mathbf{X}_{11} & \sin \mathbf{X}_{12} \\ \sin \mathbf{X}_{21} & \sin \mathbf{X}_{22} \end{bmatrix} = \begin{bmatrix} \cos \mathbf{X}_{11} d\mathbf{X}_{11} & \cos \mathbf{X}_{12} d\mathbf{X}_{12} \\ \cos \mathbf{X}_{21} d\mathbf{X}_{21} & \cos \mathbf{X}_{22} d\mathbf{X}_{22} \end{bmatrix} = \cos(\mathbf{X}) \odot d\mathbf{X}.$$

If the scalar function $f(\mathbf{X})$ is constructed through operations such as addition, subtraction, multiplication, inversion, or element-wise functions on the matrix \mathbf{X} , we can apply the above formulas to express the differential df in the form presented in Equation (2.2). To achieve this, we need some trace tricks.

2.1.2 Trace tricks

1. If $\mathbf{a} \in \mathbb{R}^{n \times 1}$, $\mathbf{B} \in \mathbb{R}^{n \times n}$,

$$\mathbf{a}^T \mathbf{B} \mathbf{a} = \text{Tr}(\mathbf{a}^T \mathbf{B} \mathbf{a}) = \text{Tr}(\mathbf{a} \mathbf{a}^T \mathbf{B}). \quad (2.3)$$

Sketch of proof: $\mathbf{a}^T \mathbf{B} \mathbf{a} = \sum_{j=1}^n a_{j1} \sum_{i=1}^n a_{i1} b_{ij} = \text{Tr}(\mathbf{a}^T \mathbf{B} \mathbf{a}) = \text{Tr}(\mathbf{a} \mathbf{a}^T \mathbf{B})$.

2. If $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{m \times m}$,

$$\text{Tr}(\mathbf{A}^T (\mathbf{B} \odot \mathbf{C})) = \text{Tr}[(\mathbf{A} \odot \mathbf{B})^T \mathbf{C}]. \quad (2.4)$$

Sketch of proof: $\text{Tr}(\mathbf{A}^T (\mathbf{B} \odot \mathbf{C})) = \sum_{i=1}^m \sum_{j=1}^m a_{ij} b_{ij} c_{ij} = \text{Tr}[(\mathbf{A} \odot \mathbf{B})^T \mathbf{C}]$.

Now the basic operation rules are prepared, to compute complex function derivative, we have one more topic to cover – composite function derivative.

2.1.3 Composite function derivative

If \mathbf{Y} is a function of \mathbf{X} and $\frac{\partial f}{\partial \mathbf{Y}}$ is known, we want to compute $\frac{\partial f}{\partial \mathbf{X}}$ using composite function derivative. In univariate calculus, we use the chain rule to compute $\frac{\partial f}{\partial x}$. But in matrix derivation, the derivative between two matrices $\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$ is undefined yet. However, we can still use the same differential operation rules to transform $d\mathbf{Y}$ into $d\mathbf{X}$. In this way, it is natural to derive the derivative $\frac{\partial f}{\partial \mathbf{X}}$. For example, if $\mathbf{Y} = \mathbf{AXB}$, we get for df ,

$$df = \text{Tr} \left[\frac{\partial f}{\partial \mathbf{Y}}^T d\mathbf{Y} \right] = \text{Tr} \left[\frac{\partial f}{\partial \mathbf{Y}}^T \mathbf{A} \text{Ad} \mathbf{X} \mathbf{B} \right] = \text{Tr} \left[\mathbf{B} \frac{\partial f}{\partial \mathbf{Y}}^T \mathbf{A} \text{Ad} \mathbf{X} \right].$$

According to the form of Equation (2.2), we obtain the derivative of f to \mathbf{X} as:

$$\frac{\partial f}{\partial \mathbf{X}} = \mathbf{A}^T \frac{\partial f}{\partial \mathbf{Y}} \mathbf{B}^T.$$

Next, we take the above methods into practice.

2.1.4 Example: logistic regression

In logistic regression, $\mathbf{y} \in \{0, 1\}^k$ is a one-hot vector acting as the label for input $\mathbf{x} \in \mathbb{R}^n$, the weight matrix is $\mathbf{W} \in \mathbb{R}^{k \times n}$. We define a probability vector $\mathbf{p} \in \mathbb{R}^k$, with p_i representing the probability of \mathbf{x} belonging to category i . The likelihood of the maximum likelihood estimation in \mathbf{p} can be expressed as:

$$\mathcal{L} = \max_{\mathbf{p}_i} \prod_{i=1}^k p_i^{y_i},$$

where y_i is the i -th element of \mathbf{y} , p_i is the i -th element of \mathbf{p} .

Next, we want to transform \prod into \sum using the log trick:

$$-\log \mathcal{L} = \min_{p_i} \left(-\sum_{i=1}^k y_i \log p_i \right),$$

where \log represents the natural logarithm.

Therefore, we define the loss function of logistic regression as:

$$l(\mathbf{x}; \mathbf{W}) = -\underbrace{\mathbf{y}^T \log \text{softmax}(\mathbf{W}\mathbf{x})}_{\mathbf{p}}. \quad (2.5)$$

To optimize l , we need to compute the derivative of l to \mathbf{W} . To simplify notations, we can view $\mathbf{W}\mathbf{x}$ as a new variable \mathbf{a} , and Equation (2.5) transforms to:

$$l(\mathbf{x}; \mathbf{W}) = -(\log \text{softmax}(\mathbf{x}^T \mathbf{W}^T)) \mathbf{y} = -(\log \text{softmax}(\mathbf{a}^T)) \mathbf{y},$$

recall that $\text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a})}{\mathbf{1}_k^T \exp(\mathbf{a})}$, where $\mathbf{1}_k$ is a k -dimensional all-ones vector, then we get for $l(\mathbf{x}; \mathbf{W})$,

$$\begin{aligned} l(\mathbf{x}; \mathbf{W}) &= -\log \left[\frac{\exp(\mathbf{a}^T)}{\exp(\mathbf{a}^T) \mathbf{1}_k} \right] \mathbf{y} \\ &= -\log [\exp(\mathbf{a}^T)] \mathbf{y} + \log [\exp(\mathbf{a}^T) \mathbf{1}_k] \mathbf{1}_k^T \mathbf{y} \quad \log(\mathbf{u}/c) = \log(\mathbf{u}) - \mathbf{1} \log(c) \\ &= -\mathbf{y}^T \mathbf{a} + \log [\exp(\mathbf{a}^T) \mathbf{1}_k]. \quad \mathbf{y}^T \mathbf{1} = 1 \end{aligned}$$

Then, we differentiate both sides of the equation on \mathbf{a} ,

$$\begin{aligned} dl &= -\mathbf{y}^T d\mathbf{a} + \frac{1}{\exp(\mathbf{a}^T) \mathbf{1}_k} [\exp(\mathbf{a}^T)] \mathbf{1}_k \\ &= -\mathbf{y}^T d\mathbf{a} + \frac{1}{\exp(\mathbf{a}^T) \mathbf{1}_k} [\exp(\mathbf{a}^T) \odot d\mathbf{a}^T] \mathbf{1}_k. \quad d\sigma(\mathbf{a}) = \sigma'(\mathbf{a}) \odot d\mathbf{a} \end{aligned}$$

According to Equation (2.2), we apply the trace operator to both sides of the equation,

$$\begin{aligned} dl &= \text{Tr} \left(-\mathbf{y}^T d\mathbf{a} + \frac{1}{\exp(\mathbf{a}^T) \mathbf{1}_k} \exp(\mathbf{a}^T) (\mathbf{d}\mathbf{a} \odot \mathbf{1}_k) \right) \\ &= \text{Tr} \left(-\mathbf{y}^T d\mathbf{a} + \frac{\exp(\mathbf{a}^T)}{\exp(\mathbf{a}^T) \mathbf{1}_k} d\mathbf{a} \right) \\ &= \text{Tr} (-[\mathbf{y}^T + \text{softmax}(\mathbf{a}^T)] d\mathbf{a}). \end{aligned}$$

Therefore,

$$\frac{\partial l}{\partial \mathbf{a}} = -\mathbf{y} + \text{softmax}(\mathbf{a}).$$

Then we apply composite function derivative rules on \mathbf{a} ,

$$dl = \text{Tr} \left(\frac{\partial l}{\partial \mathbf{a}}^T d\mathbf{a} \right) = \text{Tr} \left(\frac{\partial l}{\partial \mathbf{a}}^T d\mathbf{W}\mathbf{x} \right) = \text{Tr} \left(\mathbf{x} \frac{\partial l}{\partial \mathbf{a}}^T d\mathbf{W} \right).$$

Therefore,

$$\frac{\partial l}{\partial \mathbf{W}} = \frac{\partial l}{\partial \mathbf{a}} \mathbf{x}^T = -\mathbf{y}\mathbf{x}^T + \text{softmax}(\mathbf{a})\mathbf{x}^T.$$

2.2 Numerical analysis

2.2.1 Norm

Norm maps a vector into a scalar ‘‘magnitude’’: $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$, often written as $\|\mathbf{x}\|$. A function $\|\mathbf{x}\| : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$ is called a norm if and only if it satisfies the following conditions:

1. $\|\mathbf{x}\| = 0 \iff \mathbf{x} = 0$
2. $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$
3. $\|\mathbf{x}\| \geq 0$
4. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

A specific norm is determined with a parameter p , referred to as p -norm. If we have $\mathbf{x} \in \mathbb{R}^{n \times 1}$, the p -norm of \mathbf{x} is defined as:

$$\|\mathbf{x}\|_p^p := \sum_i^n |x_i|^p, \quad (2.6)$$

when $p = \infty$,

$$\|\mathbf{x}\|_\infty = \max_i |x_i|. \quad (2.7)$$

The ∞ -norm of a vector is the maximum absolute value of its elements.

when $p = 0$,

$$\|\mathbf{x}\|_0 = \sum_{i=1}^n \mathbf{1}\{x_i \neq 0\}, \quad (2.8)$$

where $\mathbf{1}\{\cdot\}$ is an indicator function. The 0-norm counts the number of non-zero elements in the vector. Further, we discuss **matrix norm**. We begin with the Frobenius norm, if we have $\mathbf{A} \in \mathbb{R}^{m \times n}$, the Frobenius norm of \mathbf{A} is:

$$\|\mathbf{A}\|_F^2 = \sum_{j=1}^n \sum_{i=1}^m a_{ij}^2 = \|\text{Vec}(\mathbf{A})\|^2, \quad (2.9)$$

where the $m \times n$ matrix \mathbf{A} can be viewed as the vector obtained by concatenating together the columns of \mathbf{A} , and the Frobenius norm can be viewed as applying the 2-norm on this new vector.

Next we introduce the operator norm. If \mathbf{X} and \mathbf{Y} are two vector spaces with norm $\|\mathbf{x}\|_p$ and $\|\mathbf{y}\|_q$, respectively. \mathbf{A} is the matrix that maps \mathbf{X} to \mathbf{Y} , $\mathbf{A} : \mathbf{X} \rightarrow \mathbf{Y}$. Operator norm $\|\mathbf{A}\|_{pq}$ is induced by vector norm:

$$\|\mathbf{A}\|_{pq} := \inf \{C \geq 0 \mid \|\mathbf{A}\mathbf{x}\|_q \leq C\|\mathbf{x}\|_p, \forall \mathbf{x} \in \mathbf{X}\}. \quad (2.10)$$

In this definition, $\|\mathbf{A}\|_{pq}$ is the maximum scaling factor that transforms the norm of vector \mathbf{x} in space \mathbf{X} to the norm of $\mathbf{A}\mathbf{x}$ in space \mathbf{Y} . The relative scaling effect of \mathbf{A} on \mathbf{x} is not influenced by the norm of \mathbf{x} . Therefore, if we simply consider the situation where $\|\mathbf{x}\|_p = 1$, we can get for $\|\mathbf{A}\|_{pq}$,

$$\|\mathbf{A}\|_{pq} = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_q. \quad (2.11)$$

Taking $p = q = 2$, we have the following inequality,

$$\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{x}\|_2. \quad \forall \mathbf{x} \in \mathbf{X} \quad (2.12)$$

On the unit sphere in the vector space, the norm of \mathbf{x} equals 1,

$$\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_2, \quad \forall \|\mathbf{x}\|_2 = 1. \quad \mathbf{x} \in \mathbf{X}$$

2.2.2 Conditioning

Conditioning refers to a measure of sensitivity of a function's output to input perturbations, often affecting the numerical stability and accuracy of computations. Relative condition number is defined as the maximum ratio of the relative error in the output of a function to the relative perturbation in the input. If we have an input vector $\mathbf{x} \in \mathbb{R}^n$ and a perturbation vector $\mathbf{h} \in \mathbb{R}^n$, we give the definition of condition number on function $f(\cdot)$ of \mathbf{x} as:

$$\kappa(f; \mathbf{x}, \mathbf{h}) = \frac{|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x})| / |f(\mathbf{x})|}{\|\mathbf{h}\| / \|\mathbf{x}\|} \quad (2.13)$$

$$\kappa(f) := \lim_{\epsilon \rightarrow 0} \max_{\mathbf{x}, \|\mathbf{h}\| \leq \epsilon \|\mathbf{x}\|} \kappa(f; \mathbf{x}, \mathbf{h}),$$

where the norm of \mathbf{h} is controlled by $\|\mathbf{x}\|$. Concisely, we will simply refer to the relative condition number as the condition number in the following analysis. Consider matrix transformation of \mathbf{x} , if $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{y} = f(\mathbf{x}) = \mathbf{Ax}$, then we have:

$$\begin{cases} \mathbf{y} = \mathbf{Ax} \\ \mathbf{y} + \delta\mathbf{y} = \mathbf{A}(\mathbf{x} + \delta\mathbf{x}). \end{cases}$$

Taking the norm of $\delta\mathbf{y}$, we have,

$$\|\delta\mathbf{y}\| = \|\mathbf{A}\delta\mathbf{x}\| \leq \|\mathbf{A}\| \|\delta\mathbf{x}\|.$$

We consider three cases,

– If \mathbf{A} is a square matrix and the inverse of \mathbf{A} exists, we have

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{y} \Rightarrow \|\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{y}\| \Rightarrow \frac{1}{\|\mathbf{y}\|} \leq \|\mathbf{A}^{-1}\| \frac{1}{\|\mathbf{x}\|}.$$

Multiplying this inequality with the above inequality of $\|\delta\mathbf{y}\|$, we get,

$$\frac{\|\delta\mathbf{y}\|}{\|\mathbf{y}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|}.$$

Based on Equation (2.13), we can compute the condition number of matrix \mathbf{A} as:

$$\kappa(f) = \kappa(\mathbf{A}) = \lim_{\delta\mathbf{x} \rightarrow 0} \max_{\mathbf{x}, \delta\mathbf{x}} \frac{\|\delta\mathbf{y}\| / \|\mathbf{y}\|}{\|\delta\mathbf{x}\| / \|\mathbf{x}\|} = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|. \quad (2.14)$$

– If $m < n$, consider the situation that $\mathbf{x} \perp \mathbf{A}$, which means that the n -dim vector \mathbf{x} is perpendicular to m row vectors in \mathbf{A} . In this case, $\|\mathbf{y}\| = 0$, and the condition number is:

$$\kappa(\mathbf{A}) = \lim_{\delta\mathbf{x} \rightarrow 0} \max_{\mathbf{x}, \delta\mathbf{x}} \frac{\|\delta\mathbf{y}\| / \|\mathbf{y}\|}{\|\delta\mathbf{x}\| / \|\mathbf{x}\|} = \infty. \quad (2.15)$$

– If $m > n$, then

$$\mathbf{x} = \mathbf{A}^+ \mathbf{Ax} = \mathbf{A}^+ \mathbf{y} \Rightarrow \|\mathbf{x}\| = \|\mathbf{A}^+ \mathbf{y}\| \leq \|\mathbf{A}^+\| \|\mathbf{y}\|,$$

where $\mathbf{A}^+ \mathbf{A} = \mathbf{I}$, and \mathbf{A}^+ is the pseudo-inverse matrix of \mathbf{A} .

$$\kappa(\mathbf{A}) = \lim_{\delta\mathbf{x} \rightarrow 0} \max_{\mathbf{x}, \delta\mathbf{x}} \frac{\|\delta\mathbf{y}\| / \|\mathbf{y}\|}{\|\delta\mathbf{x}\| / \|\mathbf{x}\|} = \|\mathbf{A}\| \|\mathbf{A}^+\|. \quad (2.16)$$

To compute \mathbf{A}^+ , we can use singular value decomposition (SVD) on \mathbf{A} .

Intuitively, if \mathbf{A} 's rank $r = n$ and \mathbf{A} is a square matrix, the equation $\mathbf{y} = \mathbf{Ax}$ has only one solution, and the condition number can be expressed using \mathbf{A}^{-1} . If $r < n$, we refer to the equation $\mathbf{y} = \mathbf{Ax}$ as underdetermined, there are infinite solutions for this equation. If $r = n$ and \mathbf{A} is not a square matrix, we refer to the equation $\mathbf{y} = \mathbf{Ax}$ as overdetermined, there's no solution to the equation, but we can use the least square method to compute the approximate solution.

2.3 Orthogonal matrices

Orthogonal matrices are square matrices whose rows and columns are orthonormal vectors, the transpose of an orthogonal matrix equals its inverse, we define orthogonal matrices as:

$$\mathbf{Q}^T \equiv \mathbf{Q}^{-1}. \quad (2.17)$$

We can compute the norm of an orthogonal matrix:

$$\begin{aligned} \|\mathbf{Q}\|^2 &= \max_{\|\mathbf{x}\|=1} \|\mathbf{Q}\mathbf{x}\|^2 \\ &= \max_{\|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{x} = 1, \end{aligned} \quad (2.18)$$

where $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$.

Similarly, we can derive the norm of the inverse of an orthogonal matrix:

$$\|\mathbf{Q}^{-1}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Q}^{-1} \mathbf{x}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Q}^T \mathbf{x}\| = \max_{\|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{Q} \mathbf{Q}^T \mathbf{x} = 1. \quad (2.19)$$

2.4 Singular value decomposition

SVD factorizes any matrix into three matrices consisting of two orthogonal matrices and a diagonal matrix of singular values. For matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$,

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T, \quad (2.20)$$

where \mathbf{U} and \mathbf{V} are two orthogonal matrices. The columns of \mathbf{U} are the left singular vectors of \mathbf{A} , the columns of \mathbf{V} are the right singular vectors of \mathbf{A} , and Σ is a diagonal matrix whose diagonal elements are the singular values of matrix \mathbf{A} .

The rank of \mathbf{A} satisfies $r \leq \min(n, m)$, then

$$\mathbf{A} = \mathbf{U}_{n \times r} \Sigma_{r \times r} (\mathbf{V}^T)_{r \times m} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^T, \quad (2.21)$$

where s_i is the i -th element in the diagonal of Σ , also the i -th singular value of \mathbf{A} , \mathbf{u}_i is the i -th column vector in \mathbf{U} and \mathbf{v}_i is the i -th column vector in \mathbf{V} .

This equation indicates that a matrix is the summation of the multiplication of its singular values and corresponding singular vectors. In some cases, we only need the first (maximum) k singular values and singular vectors to express \mathbf{A} , so the truncated SVD can be expressed as:

$$\tilde{\mathbf{A}} = \sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^T. \quad (2.22)$$

Next, we examine the 2-norm of \mathbf{A} from SVD perspective (Note: We omit the subscript 2 of the 2-norm in the derivation to keep the proof concise.),

$$\|\mathbf{A}\| \leq \|\mathbf{U}\| \|\Sigma\| \|\mathbf{V}^T\| = \|\Sigma\| = \sigma_{\max},$$

where $\|\mathbf{U}\| = \|\mathbf{V}\| = 1$.

Similarly, Σ can be expressed using \mathbf{A} ,

$$\Sigma = \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T \mathbf{V} = \mathbf{U}^T \mathbf{A} \mathbf{V}.$$

The norm of Σ satisfies the following inequality,

$$\|\Sigma\| \leq \|\mathbf{U}^T\| \|\mathbf{A}\| \|\mathbf{V}\| = \|\mathbf{A}\|.$$

Therefore,

$$\|\mathbf{A}\| = \|\Sigma\| = \sigma_{\max}. \quad (2.23)$$

This equation indicates that a matrix's norm equals its maximum singular value.

Now we consider the situation of $\mathbf{A}^T \mathbf{A}$, $\mathbf{A}^T \mathbf{A}$ can be expressed using the SVD form of \mathbf{A} :

$$\mathbf{A}^T \mathbf{A} = \mathbf{V} \Sigma \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T = \mathbf{V} \Sigma^2 \mathbf{V}^T. \quad (2.24)$$

This equation shows that the diagonal elements in Σ^2 are eigenvalues of $\mathbf{A}^T \mathbf{A}$.

Using SVD, the pseudo-inverse matrix of \mathbf{A} can be defined as:

$$\mathbf{A}^+ = \mathbf{V}_{m \times r} \Sigma_{r \times r}^{-1} (\mathbf{U}^T)_{r \times n}. \quad (2.25)$$

Similar to $\|\mathbf{A}\|$, we can derive the norm of \mathbf{A}^+ as:

$$\|\mathbf{A}^+\| = \frac{1}{\sigma_{\min}}. \quad (2.26)$$

2.5 Positive semi-definite

A matrix is positive semi-definite (PSD) if any quadratic form it defines yields no negative values.

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0, \quad \forall \mathbf{x} \quad (2.27)$$

A symmetric matrix is PSD if all its eigenvalues are non-negative. For a PSD matrix \mathbf{A} , it can be factorized using eigenvalue decomposition:

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{U}^T, \quad (2.28)$$

where \mathbf{U} is an orthogonal matrix, and Σ is a diagonal matrix with diagonal elements being eigenvalues of \mathbf{A} .

PSD matrices are very useful in machine learning applications. For example, in the attention mechanism, we have query matrix $\mathbf{Q} \in \mathbb{R}^{n \times d}$ and key matrix $\mathbf{K} \in \mathbb{R}^{n \times d}$, the attention score is computed using the similarity between \mathbf{Q} and \mathbf{K} , which is the inner product of \mathbf{Q} and \mathbf{K} through the exponential function, $\exp(\mathbf{Q} \mathbf{K}^T)$. If we consider a matrix \mathbf{X} composed of \mathbf{Q} and \mathbf{K} :

$$\mathbf{X} = \begin{bmatrix} \mathbf{Q} \\ \mathbf{K} \end{bmatrix}$$

Then the matrix $\exp(\mathbf{X} \mathbf{X}^T)$ is a PSD matrix with $\exp(\mathbf{Q} \mathbf{K}^T)$ as its right upper component,

$$\exp(\mathbf{X} \mathbf{X}^T) = \exp \left(\begin{bmatrix} \mathbf{Q} \mathbf{Q}^T & \mathbf{Q} \mathbf{K}^T \\ \mathbf{K} \mathbf{Q}^T & \mathbf{K} \mathbf{K}^T \end{bmatrix} \right).$$

Further, we normalize the attention score using softmax operation, we first compute the denominator:

$$D = \text{diag}(\exp(\mathbf{Q} \mathbf{K}^T) \cdot \mathbf{1}_n),$$

where $\mathbf{1}_n$ is a column vector of all ones with length n , and the diag operation refers to constructing a new matrix by placing all the elements of a given vector along its main diagonal. The attention score matrix is further computed as:

$$A = D^{-1} \exp(\mathbf{Q} \mathbf{K}^T)$$

Following the above procedure, we can see that the complexity of attention computation is quadratic with sequence length n .

2.6 Revisit linear regression

Recall that the optimization objective of a linear regression model can be described as the equation below:

$$\beta^* = \arg \min_{\beta} \langle \mathbf{X}\beta - \mathbf{Y}, \mathbf{X}\beta - \mathbf{Y} \rangle. \quad (2.29)$$

We make the inner product term as a function $f(\beta)$, then take the derivative of the square loss using matrix derivative rules,

$$\frac{\partial f}{\partial \beta} = 0 \Rightarrow \mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \mathbf{Y}. \quad (2.30)$$

To examine the numerical stability of $\hat{\beta}$, we compute the condition number of $\mathbf{X}^T \mathbf{X}$ based on Equation (2.14):

$$\kappa(\mathbf{X}^T \mathbf{X}) = \|\mathbf{X}\|^2 \|\mathbf{X}^+\|^2. \quad (2.31)$$

Using QR factorization $\mathbf{X} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is an orthogonal matrix and \mathbf{R} is an upper triangular matrix. Then we have,

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \hat{\beta} &= \mathbf{X}^T \mathbf{Y} \Rightarrow \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R} \hat{\beta} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Y} \\ &\Rightarrow \mathbf{R}^T \mathbf{R} \hat{\beta} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Y} \\ &\Rightarrow \mathbf{R} \hat{\beta} = \mathbf{Q}^T \mathbf{Y} \end{aligned}$$

The condition number for this linear system is

$$\kappa(\mathbf{R}) = \kappa(\mathbf{X}) = \|\mathbf{X}\| \|\mathbf{X}^+\|. \quad (2.32)$$

From the above computation, we can know that using QR factorization significantly reduces condition number, a trick we can use for estimation.

Further, we check the variance for $\hat{\beta}$ from the perspective of conditioning. We represent the condition number using matrix singular value. According to Equation (2.23) and Equation (2.26), the condition number $\kappa(\mathbf{X}^T \mathbf{X})$ can be further expressed as

$$\kappa(\mathbf{X}^T \mathbf{X}) = \kappa(\mathbf{V}\Sigma^2\mathbf{V}^T) = \kappa(\Sigma^2) = \frac{\sigma_{\max}^2}{\sigma_{\min}^2}. \quad (2.33)$$

Revisit the variance of $\hat{\beta}$,

$$\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \quad (2.34)$$

If we fix the 2-norm of \mathbf{X} as 1, then the maximum singular value σ_{\max} equals to 1. The condition number of $\mathbf{X}^T \mathbf{X}$ can be re-written as:

$$\kappa(\mathbf{X}^T \mathbf{X}) = \frac{1}{\sigma_{\min}^2}.$$

Taking the norm of variance on $\hat{\beta}$, we have:

$$\|\text{Var}(\hat{\beta})\| = \|\sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\| = \frac{\sigma^2}{\sigma_{\min}^2}.$$

In this case, when the smallest singular value σ_{\min} of \mathbf{X} is small, it leads to a large condition number and high variance. This situation indicates the instability of estimation when the collinearity among variables is high.