

无约束的矩阵低秩近似 伪逆与SVD奇异值分解 LoRA微调和MLA的数学基础

东川路第一可爱猫猫虫



主要内容

- 矩阵的范数
- Frobenius范数
- 伪逆（广义逆）
- 最优化求伪逆
- SVD分解求低秩近似
- Eckart-Young-Mirsky 定理

矩阵的范数

- 范数是向量或矩阵的一种度量

范数是将矩阵映射为非负实数的函数

若 $N(A) = ||A||$ 满足正定性，齐次性，三角不等式，则称
 $N(A) = ||A||$ 是矩阵 A 的范数

- Frobenius 范数

$$||M||_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m M_{i,j}^2}$$

把矩阵的所有元素的平方加起来，开根号

简称 F 范数

正交变换不改变 F 范数

- $\|A\|_F^2 = \text{tr}(A^T A)$

proof: $A^T A$ 展开, 其主对角线上元素为第*i*行元素平方和, 再对主对角线上元素求和即得到 $\|A\|_F^2$

- 对任意矩阵 ($A \in R^{m \times n}$), 左乘或右乘正交阵均称为正交变换
- 以左乘正交阵为例来证明

$$\|QA\|_F^2 = \text{tr}[(QA)^T (QA)] = \text{tr}[A^T Q^T QA] = \text{tr}[A^T A] = \|A\|_F^2$$

右乘正交阵的情况同理可证

广义逆 (伪逆)

- 逆矩阵

若矩阵 $AB=M$, 且 A 为可逆的方阵, 则有 $B=A^{-1}M$ 称为 A 的逆
那如果 A 不可逆呢? 甚至 A 不是方阵呢?

- 广义逆

1920年Moore提出广义逆

1954年Penrose独立地给出了更完备的定义

$$AXA = A,$$

这四个等式存在唯一解, 这个解 X 就是矩阵 A 的伪逆

$$XAX = X,$$

《A GENERALIZED INVERSE FOR MATRICES》

$$(AX)^* = AX,$$

$$(XA)^* = XA,$$

另一种方式求伪逆

- 在已知矩阵A和M的情况下，求出一个B使得AB与M间误差最小

$$A \text{的伪逆} = \operatorname{argmin}_B \|AB - M\|_F^2$$

其中 $A \in R^{n \times r}, B \in R^{r \times m}, M \in R^{n \times m}$

上式的运算结果是一个使得 $L = \|AB - M\|_F^2$ 最小的矩阵B

为使 $L = \|AB - M\|_F^2$ 最小，求L对矩阵B的导数

- 标量函数f的矩阵X的导数

$$\frac{\partial f}{\partial X} = \left[\frac{\partial f}{\partial x_{i,j}} \right]_{m \times n}$$

求 $L = \left\| AB - M \right\|_F^2$ 对矩阵B的导数

- 链式法则

令 $E = AB - M$

$$\frac{\partial L}{\partial B_{i,j}} = \sum_{k,l} E_{k,l} \frac{\partial L}{\partial E_{k,l}} \frac{\partial E_{k,l}}{\partial B_{i,j}}$$

这里 $E_{k,l}$ 代表矩阵E在(k,l)位置的元素， $B_{i,j}$ 代表矩阵B在(ki,j)位置的元素

接下来分别求 $\frac{\partial L}{\partial E_{k,l}}$ 和 $\frac{\partial E_{k,l}}{\partial B_{i,j}}$

求第一部分 $\frac{\partial L}{E_{k,l}}$

- 由 $L = \left\| AB - M \right\|_F^2$ 以及 F 范数的定义

$$L = \left\| E \right\|_F^2 = \sum_{i,j} E_{i,j}^2$$

多元函数对一个元求偏导，只有 $(l,j) = (k,l)$ 时， $\frac{\partial L}{E_{k,l}}$ 才不为 0

也就是说， L 关于 $E_{k,l}$ 的导数就是 $E_{k,l}^2$ 关于 $E_{k,l}$ 的导数

$$\frac{\partial L}{\partial E_{k,l}} = 2E_{k,l}, \quad \frac{\partial L}{\partial E} = 2E$$

求第二部分 $\frac{\partial E_{k,l}}{\partial B_{i,j}}$

- $E = AB - M$ 而 M 是常数矩阵，与 B 无关 因此 $\frac{\partial E}{\partial B} = \frac{\partial AB}{\partial B}$

- $A_{n \times r} B_{r \times m}$ 在 (k,l) 位置上的元素为？

这个元素实际上是 A 的第 k 行向量与 B 的第 l 列向量做内积

$$(A_{k,1}, \dots, A_{k,r})(B_{1,l}, \dots, B_{r,l}) = \sum_m A_{k,m} B_{m,l}$$

- 由此 E 在 (k,l) 位置的元素为

$$\sum_m A_{k,m} B_{m,l} - M_{k,l}$$

若 $l \neq j$, 则 $\frac{\partial E_{k,l}}{\partial B_{i,j}} = 0$; 若 $l = j$, 则 $\frac{\partial E_{k,l}}{\partial B_{i,j}} = A_{k,i}$

$$\frac{\partial E_{k,l}}{\partial B_{i,j}} = A_{k,i} \delta_{l,j} \quad \text{其中 } \delta_{l,j} \text{ 用于判定 } l \text{ 与 } j \text{ 是否相等}$$

$$\text{求} \frac{\partial L}{\partial B_{i,j}} = \sum_{k,l} \frac{\partial L}{E_{k,l}} \frac{\partial E_{k,l}}{\partial B_{i,j}}$$

- $\frac{\partial L}{\partial B_{i,j}} = 2 \sum_{k,l} E_{k,l} A_{k,i} \delta_{l,j} = 2 \sum_k E_{k,j} A_{k,i}$
- 注意到 $\sum_k E_{k,j} A_{k,i}$ 是 $A^T E$ 在(i,j)的元素

因此有 $\frac{\partial L}{\partial B_{i,j}} = 2(A^T E)_{i,j}$

$$\frac{\partial L}{\partial B} = 2(A^T E) = 2A^T(AB - M)$$

同理, $\frac{\partial L}{\partial A} = 2(AB - M)B^T$

要求min, 令偏导=0

- $\frac{\partial L}{\partial B} = 2A^T(AB - M) = 0, \quad \frac{\partial L}{\partial A} = 2(AB - M)B^T$

有 $A^T AB = A^T M, \quad ABB^T = MB^T$

只要 $A^T A$ 可逆, 我们就求出了 $B = (A^T A)^{-1} A^T M$

这就是我们求得的 $\operatorname{argmin}_B \|AB - M\|_F^2$, A 的伪逆

以上就是给定矩阵 M 和 A 时, 优化 $\|AB - M\|_F^2$ 的最优解

A和B均未知的话， 该怎么办

- 刚才我们在已知矩阵A和M的情况下， 求出了 $\operatorname{argmin}_B \|\mathbf{AB} - \mathbf{M}\|_F^2$
- 现在考虑A、B未知时候的最优解

即求 $\operatorname{argmin}_{\mathbf{A}, \mathbf{B}} \|\mathbf{AB} - \mathbf{M}\|_F^2$

其中， $\mathbf{A} \in \mathbb{R}^{n \times r}, \mathbf{B} \in \mathbb{R}^{r \times m}, \mathbf{M} \in \mathbb{R}^{n \times m}, r < \min(n, m)$

- SVD奇异值分解

奇异值分解在实矩阵域的定理：

对于任意矩阵 $\mathbf{M} \in \mathbb{R}^{n \times m}$ ， 都能找到 $\mathbf{M} = \mathbf{U} \Sigma \mathbf{V}^T$

其中 \mathbf{U} 、 \mathbf{V} 为 $n \times n$ 、 $m \times m$ 正交阵， Σ 是非负对角阵
对角线元素称为奇异值， 默认从大到小排序

用SVD求解 $argmin_{A,B} ||AB - M||_F^2$

- 将M分解为 $U\Sigma V^T$

$$\begin{aligned} \text{则} ||AB - M||_F^2 &= ||UU^T ABV V^T - U\Sigma V^T||_F^2 \\ &= ||U(U^T ABV - \Sigma)V^T||_F^2 \\ &= ||U^T ABV - \Sigma||_F^2 \end{aligned}$$

问题被化简了

- 低秩近似

低秩近似的目地是求“M的最优r秩近似”

即求秩不超过r的矩阵X，使得 $\|X - M\|_F^2$ 最小

矩阵的秩分解性质：秩为r的矩阵可分解为“ $m \times r$ 列满秩矩阵”与“ $r \times n$ 行满秩矩阵”的乘积

因此问题可以转化为 $\operatorname{argmin}_{A,B} \|AB - M\|_F^2$

由刚才的SVD分解，我们进一步把问题转化为了：

$$\operatorname{argmin}_{A,B} \|U^T ABV - \Sigma\|_F^2$$

问题的进一步化简

- 令 $Y = U^T A B V$
- 我们注意到

Y 的秩与 AB 相同 (正交变换不改变矩阵的秩) $\leq r$

任意秩不超过 r 的矩阵 Y , 都能找到 A 和 B , 这是因为 U 和 V 都是可逆的, $AB = U Y V^T$, 而 UYV^T 秩 $\leq r$, 又可以分解为 A 和 B

- 由此, 问题被简化成了非负对角阵 Σ 的最优 r 秩近似

非负对角阵 Σ 的最优 r 秩近似

- 我们要求出一个秩 $\leq r$ 的矩阵 Σ_r , 以最小化 $\|\Sigma_r - \Sigma\|_F^2$
其中 Σ 为已知的非负对角阵
- 当然我们考虑的是低秩近似, 如果 Σ 秩 $\leq r$, 那就取 $\Sigma_r = \Sigma$, 这样 $\|\Sigma_r - \Sigma\|_F^2$ 直接为0了, 结论很显然, 没有讨论价值
- 所以我们考虑 Σ 秩 $> r$ 的情况

为达到目标, Σ_r 的非对角线元素均需为0, 消除非对角线误差
对角线上, 保留前 r 个最大的奇异值, 其余为0
- 结论

非负对角阵的最优 r 秩近似就是只保留对角线最大的 r 个元素的矩阵

基于 SVD 的Eckart-Young-Mirsky 定理

- 如果 $M \in R^{m \times n}$ 的SVD奇异值分解为 $U\Sigma V^T$, 那么 M 的最优 r 秩近似为 $U_{[:n,:r]}\Sigma_{[:r,:r]}V_{[:m,:r]}^T$
 - $U_{[:n,:r]}$ 取正交矩阵 U 的前 r 列 (保留前 r 个“左特征方向”)
 - $\Sigma_{[:r,:r]}$ 取对角阵 Σ 的前 r 阶对角子阵, 保留前 r 个最大的奇异值
 - $V_{[:m,:r]}^T$ 取正交矩阵 V^T 的前 r 列 (保留前 r 个“右特征方向”)
- 普遍意义在于:
 - 任何矩阵 (无论是否对角、是否方阵) 的最优 r 秩近似, 都可以通过其 SVD 分解, “截取前 r 个最大奇异值及对应方向” 得到