

符号求导原理

梯度 链式法则

符号求导原理

问题原型

求导原理分析

变分原理与表达式 \boldsymbol{h}_{ij} 的推导

多输入单输出时的梯度算法

多输出单输入时的梯度算法

多输出多输入时的梯度算法

离散变量求导法

问题原型

假设存在如下的运算：

$$\boldsymbol{Y}_1, \cdots, \boldsymbol{Y}_m = f(\boldsymbol{X}_1, \cdots, \boldsymbol{X}_n)$$

其中 $\boldsymbol{Y}_1, \cdots, \boldsymbol{Y}_m, \boldsymbol{X}_1, \cdots, \boldsymbol{X}_n$ 均为张量,且 f 不一定可以解耦成函数列. 现在对于任意作用于 $\boldsymbol{Y}_1, \cdots, \boldsymbol{Y}_m$ 的标量函数 g ：

$$z = g(\boldsymbol{Y}_1, \cdots, \boldsymbol{Y}_m)$$

若已知 $\frac{\partial z}{\partial \boldsymbol{Y}_1}, \cdots, \frac{\partial z}{\partial \boldsymbol{Y}_m}, \boldsymbol{Y}_1, \cdots, \boldsymbol{Y}_m, \boldsymbol{X}_1, \cdots, \boldsymbol{X}_n$ 问是否可以计算出 $\frac{\partial z}{\partial \boldsymbol{X}_1}, \cdots, \frac{\partial z}{\partial \boldsymbol{X}_n}$ ？
其实问题是肯定的.

求导原理分析

设 $\boldsymbol{Y}, \boldsymbol{X}$ 均为张量,那么 $\frac{\partial \boldsymbol{Y}}{\partial \boldsymbol{X}}$ 依 \boldsymbol{Y} 元素梯度的结果可以排列成一个新的张量,其shape为 $\boldsymbol{Y}, \boldsymbol{X}$ 二者shape的叠加.则普通的链式法则可以表示成张量积的形式：

$$\frac{\partial z}{\partial \boldsymbol{X}_i} = \sum_j \frac{\partial z}{\partial \boldsymbol{Y}_j} \frac{\partial \boldsymbol{Y}_j}{\partial \boldsymbol{X}_i}, \quad i = 1, \cdots, n$$

由于 $\frac{\partial \boldsymbol{Y}_j}{\partial \boldsymbol{X}_i}$ 表达式肯定是与 g 无关的,所以(1)可以被表达为：

$$\frac{\partial z}{\partial \boldsymbol{X}_i} = \sum_j \frac{\partial z}{\partial \boldsymbol{Y}_j} \boldsymbol{H}_{ij}(\boldsymbol{X}_1, \cdots, \boldsymbol{X}_n) = \sum_j \boldsymbol{h}_{ij}(\frac{\partial z}{\partial \boldsymbol{Y}_j}, \boldsymbol{X}_1, \cdots, \boldsymbol{X}_n), \quad i = 1, \cdots, n$$

其中 $\boldsymbol{H}_{ij}, \boldsymbol{h}_{ij}$ 均为输出为张量的函数. \boldsymbol{H}_{ij} 的shape为 \boldsymbol{Y}_j 和 \boldsymbol{X}_i shape的叠加, \boldsymbol{h}_{ij} 的shape为 \boldsymbol{X}_i shape的一致.从(2)可以看出 \boldsymbol{H}_{ij} 的表达式与 g 无关,进而 \boldsymbol{h}_{ij} 的表达式也与 g 无关.

所以若已知 $\frac{\partial z}{\partial \boldsymbol{Y}_1}, \cdots, \frac{\partial z}{\partial \boldsymbol{Y}_m}, \boldsymbol{Y}_1, \cdots, \boldsymbol{Y}_m, \boldsymbol{X}_1, \cdots, \boldsymbol{X}_n$,是可以通过推导表达式 \boldsymbol{h}_{ij} 计算出 $\frac{\partial z}{\partial \boldsymbol{X}_1}, \cdots, \frac{\partial z}{\partial \boldsymbol{X}_n}$ 的.

变分原理与表达式 \boldsymbol{h}_{ij} 的推导

我们参考PDE里的变分原理,来实现对表达式 \boldsymbol{h}_{ij} 的推导.

多输入单输出时的梯度算法

首先考虑 $m = 1$ 的情况.不失一般性,记 \boldsymbol{Y}_1 是一个含有三个上标的张量,构造一个下标与 \boldsymbol{Y}_1 上标完全一致的任意张量 $\Phi^{(1)}$.通过Einstein求和,这样就可以构造出不可数个 g ：

$$g(\boldsymbol{Y}_1) = \Phi_{ijk}^{(1)} \boldsymbol{Y}_1^{ijk}$$

我们这样做的目的是为了使得 $\frac{\partial z}{\partial \boldsymbol{Y}_1}$ 正好成为一个形状指定但是数值任意的张量：

$$\frac{\partial z}{\partial \boldsymbol{Y}_1^{ijk}} = \Phi_{ijk}^{(1)} \Rightarrow \frac{\partial z}{\partial \boldsymbol{Y}_1} = \Phi^{(1)}$$

通过将(2.2)带入1.2即可得到 $\boldsymbol{h}_{i1}, i = 1, \cdots, n$ 的表达式.

以下以 f 是矩阵乘法为例：

$$\begin{cases} Y_1 = f(X_1, X_2) = X_1 X_2 \\ z = \sum_{i,j} \Phi_{ij}^{(1)} \sum_k X_1^{ik} X_2^{kj} \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial Y_1^{ij}} = \Phi_{ij}^{(1)} \\ \frac{\partial z}{\partial X_1^{ik}} = \sum_j \Phi_{ij}^{(1)} X_2^{kj} \\ \frac{\partial z}{\partial X_2^{kj}} = \sum_i \Phi_{ij}^{(1)} X_1^{ik} \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial X_1^{ik}} = \sum_j \frac{\partial z}{\partial Y_1^{ij}} X_2^{kj} \\ \frac{\partial z}{\partial X_2^{kj}} = \sum_i \frac{\partial z}{\partial Y_1^{ij}} X_1^{ik} \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial X_1} = \frac{\partial z}{\partial Y_1} X_2^T \\ \frac{\partial z}{\partial X_2} = X_1^T \frac{\partial z}{\partial Y_1} \end{cases}$$

多输出单输入时的梯度算法

考虑 $n = 1$ 的情况,一般来说这种情况下 h_{m1} 的表达式是很难推导出来的,但是如果 [假设 \$f\$ 可逆](#),则问题将变得简单,先根据[多输入单输出时的梯度算法](#)可以得到如下的表达式组:

$$\begin{cases} \frac{\partial z}{\partial Y_1} = \hat{h}_{11}(\frac{\partial z}{\partial X_1}, Y_1, \dots, Y_m) \\ \vdots \\ \frac{\partial z}{\partial Y_m} = \hat{h}_{m1}(\frac{\partial z}{\partial X_1}, Y_1, \dots, Y_m) \end{cases}$$

通过对 [\(2.2.1\)](#) 式进行变形即可得到相应的 [\(1.2\)](#).

以下以 f 是矩阵的最大主元LU分解为例:

$$\begin{cases} LU = Pf(X) \\ z = \sum_{i,j} \Phi_{ij} X_{ij} \\ = \sum_{i,j} \Phi_{ij} \sum_k P_{ik}^T L_{kl} U_{lj} \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial X} = \Phi \\ \frac{\partial z}{\partial P} = LU\Phi^T \\ \frac{\partial z}{\partial L} = P\Phi U^T \\ \frac{\partial z}{\partial U} = L^T P\Phi \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial X} = (\frac{\partial z}{\partial P})^T L^{-T} U^{-T} \\ \frac{\partial z}{\partial X} = P^T \frac{\partial z}{\partial L} U^{-T} \\ \frac{\partial z}{\partial X} = P^T L^{-T} \frac{\partial z}{\partial U} \end{cases}$$

本例说明了一个问题,也就是针对多输出的情况, [\(1.2\)](#) 可以有多种不同的表达式且不一定需要所有的上一级梯度.

多输出多输入时的梯度算法

经过测试一般来说在这种情况下并没有什么非常简单的办法可以导出 [\(1.2\)](#),若存在 f_X 和可逆的 f_Y 使得:

$$W = f_Y(Y_1, \dots, Y_m) = f_X(X_1, \dots, X_n)$$

则可以综合[多输入单输出时的梯度算法](#)和[多输出单输入时的梯度算法](#)仍然可以导出 [\(1.2\)](#).

以下以 f 是对矩阵乘法结果的最大主元LU分解为例:

$$\begin{cases} LU = Pf(AB) \\ z = \sum_{i,j} \Phi_{ij} \sum_k A_{ik} B_{kj} \\ = \sum_{i,j} \Phi_{ij} \sum_k P_{ik}^T L_{kl} U_{lj} \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial A} = \frac{\partial z}{\partial(AB)} B^T \\ \frac{\partial z}{\partial B} = A^T \frac{\partial z}{\partial(AB)} \\ \frac{\partial z}{\partial(AB)} = (\frac{\partial z}{\partial P})^T L^{-T} U^{-T} \\ \frac{\partial z}{\partial(AB)} = P^T \frac{\partial z}{\partial L} U^{-T} \\ \frac{\partial z}{\partial(AB)} = P^T L^{-T} \frac{\partial z}{\partial U} \end{cases} \Rightarrow \begin{cases} \frac{\partial z}{\partial A} = (\frac{\partial z}{\partial P})^T L^{-T} U^{-T} B^T \\ \frac{\partial z}{\partial A} = P^T \frac{\partial z}{\partial L} U^{-T} B^T \\ \frac{\partial z}{\partial A} = P^T L^{-T} \frac{\partial z}{\partial U} B^T \\ \frac{\partial z}{\partial B} = A^T (\frac{\partial z}{\partial P})^T L^{-T} U^{-T} \\ \frac{\partial z}{\partial B} = A^T P^T \frac{\partial z}{\partial L} U^{-T} \\ \frac{\partial z}{\partial B} = A^T P^T L^{-T} \frac{\partial z}{\partial U} \end{cases}$$

离散变量求导法