

# Week 1 ML 学习总结 机器学习的环境与数学基础

## 机器学习环境搭建

笔者采用的环境/框架是 Python CUDA Torch

### 具体步骤

1. 安装Anaconda，部署Python, Jupyter Notebook等的环境。
2. 在Conda的Console中使用 `conda install` 命令安装 `pytorch` 的GPU版本。

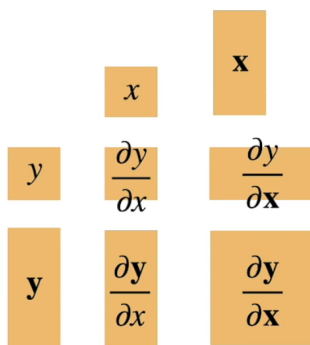
```
# for windows  
conda install pytorch torchvision torchaudio cudatoolkit=11.3 -c pytorch
```

完成。

## 定义

### 梯度

梯度是对导数的扩充。



标量函数对向量的梯度

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \frac{\partial y}{\partial \mathbf{x}} = \left[ \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n} \right] \quad (1)$$

梯度指向了值变化最大的方向。

实例

$$\frac{\partial (x_1^2 + 2x_2^2)}{\partial \mathbf{x}} = [2x_1, 4x_2] \quad (2)$$

## 标量函数对向量的梯度

### 例子

$$\begin{array}{c|cccc} y & a & au & \text{sum}(\mathbf{x}) & \|\mathbf{x}\|^2 \\ \hline \frac{\partial y}{\partial \mathbf{x}} & \mathbf{0}^T & a \frac{\partial u}{\partial \mathbf{x}} & \mathbf{1}^T & 2\mathbf{x}^T \end{array} \quad (3)$$

$$\begin{array}{c|ccc} y & u+v & uv & \langle \mathbf{u}, \mathbf{v} \rangle \\ \hline \frac{\partial y}{\partial \mathbf{x}} & \frac{\partial u}{\partial \mathbf{x}} + \frac{\partial v}{\partial \mathbf{x}} & \frac{\partial u}{\partial \mathbf{x}} v + \frac{\partial v}{\partial \mathbf{x}} u & \mathbf{u}^T \frac{\partial \mathbf{v}}{\partial \mathbf{x}} + \mathbf{v}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \end{array} \quad (4)$$

### ✎ 一些证明

$$\begin{aligned} \frac{\partial \|\mathbf{x}\|^2}{\partial \mathbf{x}} &= \frac{\partial x_1^2 + x_2^2 + \cdots + x_n^2}{\partial \mathbf{x}} \\ &= [2x_1 \quad 2x_2 \quad \cdots \quad 2x_n] \\ &= 2\mathbf{x}^T \end{aligned} \quad (5)$$

$$\begin{aligned} \frac{\partial \langle \mathbf{u}, \mathbf{v} \rangle}{\partial \mathbf{x}} &= \frac{\partial \mathbf{u}_1 \mathbf{v}_1 + \mathbf{u}_2 \mathbf{v}_2 + \cdots + \mathbf{u}_n \mathbf{v}_n}{\partial \mathbf{x}} \\ &= \mathbf{u}_1 \frac{\partial \mathbf{v}_1}{\partial \mathbf{x}} + \mathbf{v}_1 \frac{\partial \mathbf{u}_1}{\partial \mathbf{x}} + \mathbf{u}_2 \frac{\partial \mathbf{v}_2}{\partial \mathbf{x}} + \mathbf{v}_2 \frac{\partial \mathbf{u}_2}{\partial \mathbf{x}} + \cdots + \mathbf{u}_n \frac{\partial \mathbf{v}_n}{\partial \mathbf{x}} + \mathbf{v}_1 \frac{\partial \mathbf{u}_n}{\partial \mathbf{x}} \\ &= \mathbf{u}^T \frac{\partial \mathbf{v}}{\partial \mathbf{x}} + \mathbf{v}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \end{aligned} \quad (6)$$

## 向量函数对标量的梯度

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad \frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix} \quad (7)$$

## 向量函数对向量的梯度

$$\begin{aligned} \mathbf{x} &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \\ \frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \frac{\partial y_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \end{aligned} \quad (8)$$

### 例子

$$\begin{array}{c|cccc} y & a & \mathbf{x} & \mathbf{A}\mathbf{x} & \mathbf{x}^T \mathbf{A} \\ \hline \frac{\partial y}{\partial \mathbf{x}} & \mathbf{0} & \mathbf{I} & \mathbf{A} & \mathbf{A}^T \end{array} \quad (9)$$

$$\begin{array}{c|ccc} y & a\mathbf{u} & \mathbf{A}\mathbf{u} & \mathbf{u} + \mathbf{v} \\ \hline \frac{\partial \mathbf{y}}{\partial \mathbf{x}} & a \frac{\partial \mathbf{u}}{\partial \mathbf{x}} & \mathbf{A} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} & \frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \frac{\partial \mathbf{v}}{\partial \mathbf{x}} \end{array} \quad (10)$$

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial x_1}{\partial \mathbf{x}} \\ \frac{\partial x_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial x_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial x_1}{\partial x_1} & \frac{\partial x_1}{\partial x_2} & \cdots & \frac{\partial x_1}{\partial x_n} \\ \frac{\partial x_2}{\partial x_1} & \frac{\partial x_2}{\partial x_2} & \cdots & \frac{\partial x_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_n}{\partial x_1} & \frac{\partial x_n}{\partial x_2} & \cdots & \frac{\partial x_n}{\partial x_n} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = I_{n \times n} \quad (11)$$

$$\begin{aligned} \frac{\partial \mathbf{Ax}}{\partial \mathbf{x}} &= \frac{\partial \begin{bmatrix} a_{1,1}x_1 + \cdots + a_{1,n}x_n \\ a_{2,1}x_1 + \cdots + a_{2,n}x_n \\ \vdots \\ a_{m,1}x_1 + \cdots + a_{m,n}x_n \end{bmatrix}}{\partial \mathbf{x}} \\ &= \begin{bmatrix} \frac{\partial a_{1,1}x_1 + \cdots + a_{1,n}x_n}{\partial x_1} & \frac{\partial a_{1,1}x_1 + \cdots + a_{1,n}x_n}{\partial x_2} & \cdots & \frac{\partial a_{1,1}x_1 + \cdots + a_{1,n}x_n}{\partial x_n} \\ \frac{\partial a_{2,1}x_1 + \cdots + a_{2,n}x_n}{\partial x_1} & \frac{\partial a_{2,1}x_1 + \cdots + a_{2,n}x_n}{\partial x_2} & \cdots & \frac{\partial a_{2,1}x_1 + \cdots + a_{2,n}x_n}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial a_{m,1}x_1 + \cdots + a_{m,n}x_n}{\partial x_1} & \frac{\partial a_{m,1}x_1 + \cdots + a_{m,n}x_n}{\partial x_2} & \cdots & \frac{\partial a_{m,1}x_1 + \cdots + a_{m,n}x_n}{\partial x_n} \end{bmatrix} \\ &= \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \\ &= A \end{aligned} \quad (12)$$

$$\begin{aligned} \frac{\partial \mathbf{x}^T \mathbf{A}}{\partial \mathbf{x}} &= \frac{\partial \begin{bmatrix} a_{1,1}x_1 + a_{2,1}x_2 + \cdots + a_{n,1}x_n \\ a_{1,2}x_1 + a_{2,2}x_2 + \cdots + a_{n,2}x_n \\ \vdots \\ a_{1,m}x_1 + a_{2,m}x_2 + \cdots + a_{n,m}x_n \end{bmatrix}^T}{\partial \mathbf{x}} \\ &= \begin{bmatrix} \frac{a_{1,1}x_1 + a_{2,1}x_2 + \cdots + a_{n,1}x_n}{\partial x_1} & \frac{a_{1,1}x_1 + a_{2,1}x_2 + \cdots + a_{n,1}x_n}{\partial x_2} & \cdots & \frac{a_{1,1}x_1 + a_{2,1}x_2 + \cdots + a_{n,1}x_n}{\partial x_n} \\ \frac{a_{1,2}x_1 + a_{2,2}x_2 + \cdots + a_{n,2}x_n}{\partial x_1} & \frac{a_{1,2}x_1 + a_{2,2}x_2 + \cdots + a_{n,2}x_n}{\partial x_2} & \cdots & \frac{a_{1,2}x_1 + a_{2,2}x_2 + \cdots + a_{n,2}x_n}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{a_{1,m}x_1 + a_{2,m}x_2 + \cdots + a_{n,m}x_n}{\partial x_1} & \frac{a_{1,m}x_1 + a_{2,m}x_2 + \cdots + a_{n,m}x_n}{\partial x_2} & \cdots & \frac{a_{1,m}x_1 + a_{2,m}x_2 + \cdots + a_{n,m}x_n}{\partial x_n} \end{bmatrix} \\ &= \begin{bmatrix} a_{1,1} & a_{2,1} & \cdots & a_{n,1} \\ a_{1,2} & a_{2,2} & \cdots & a_{n,2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1,m} & a_{2,m} & \cdots & a_{n,m} \end{bmatrix} \\ &= A^T \end{aligned} \quad (13)$$

### 标量函数对矩阵的梯度

$m$ 维行向量函数  $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \cdots, f_m(\mathbf{x})]$  相对于  $n$  维实向量  $\mathbf{x}$  的梯度为一  $n \times m$  矩阵, 定义为

$$\begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_1} \\ \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n} & \frac{\partial f_2(\mathbf{x})}{\partial x_n} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix}$$

$$\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) \stackrel{\text{def}}{=} \begin{bmatrix} \overline{\frac{\partial}{\partial x_2}} & \overline{\frac{\partial}{\partial x_2}} & \cdots & \overline{\frac{\partial}{\partial x_2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n} & \frac{\partial f_2(\mathbf{x})}{\partial x_n} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \quad (14)$$

(关于为什么列向量变为了行向量：行向量和列向量乘积是标量)

## 求导

### 链式法则

#### 标量链式法则

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x} \quad (15)$$

#### 向量链式法则

$$\begin{aligned} \frac{\partial y}{\partial \mathbf{x}} &= \frac{\partial y}{\partial u} \frac{\partial u}{\partial \mathbf{x}} \\ (1, n) \quad (1)(1, n) \\ \frac{\partial y}{\partial \mathbf{x}} &= \frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \\ (1, n) \quad (1, k)(k, n) \\ \frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= \frac{\partial \mathbf{y}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \\ (m, n) \quad (m, k)(k, n) \end{aligned} \quad (16)$$

#### 例子

$$\begin{aligned} z &= (\langle \mathbf{x}, \mathbf{w} \rangle - y)^2 \\ a &= \langle \mathbf{x}, \mathbf{w} \rangle \\ b &= a - y \\ z &= b^2 \end{aligned} \quad (17)$$

$$\begin{aligned} \frac{\partial z}{\partial w} &= \frac{\partial z}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial w} \\ &= 2b \cdot 1 \cdot \mathbf{x}^T \quad (u, v \text{ 无关}) \\ &= (2\langle \mathbf{x}, \mathbf{w} \rangle - y) \mathbf{x}^T \\ z &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \end{aligned}$$

$$\begin{aligned} a &= \mathbf{X}\mathbf{w} \\ b &= a - \mathbf{y} \\ z &= \|\mathbf{b}\|^2 \end{aligned} \quad (18)$$

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{w}} &= \frac{\partial z}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{w}} \\ &= 2\mathbf{b}^T \cdot 1 \cdot \mathbf{X} \\ &= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T \mathbf{X} \end{aligned}$$

# 动手做！

[自动求导实验](#)

## 线性回归模型：以预测房价为例

---

### 概述

$$\begin{aligned}\text{Input: } X &= [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \\ \text{Output: } y &= [y_1, y_2, \dots, y_n]^T\end{aligned}\tag{19}$$

输入关于房产信息的向量 $x_i$ ，输出房价 $y_i$

假设对于房价的影响由三个因素确定： $x_1, x_2, x_3$

假设成交价是关键因素的加权和 $y = w_1x_1 + w_2x_2 + w_3x_3 + b$

### 推广

广泛的，可以如此表示线性模型：

$$y = \sum_{i=1}^n w_i x_i + b\tag{20}$$

也可以以向量形式表示为：

$$y = \langle \mathbf{w}, \mathbf{x} \rangle + b\tag{21}$$

### 衡量与评估质量

#### 损失函数

$$\ell(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2\tag{22}$$

(平方损失)

### 定义

#### 训练损失

$$\ell(x, y, w, b) = \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \mathbf{w} \rangle - b)^2 = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w} - b\|^2\tag{23}$$

#### 最小化损失来学习参数

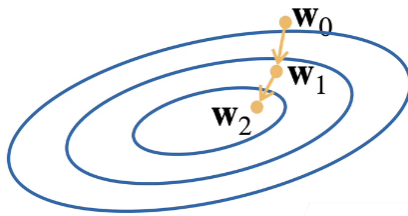
$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \ell(\mathbf{X}, \mathbf{y}, \mathbf{w}, b)\tag{24}$$

## 梯度下降

---

1. 选取初始值 $\mathbf{w}_0$
2. 迭代 $t = 1, 2, 3 \dots$

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \frac{\partial \ell}{\partial \mathbf{w}_{t-1}}\tag{25}$$



每次向着梯度的反方向前进，会最大的减少损失函数值。

$\eta$ :学习率 步长的**超参数**

\*超参数：在开始学习过程之前设置值的参数

学习率不应该过小，否则梯度下降过慢；学习率过大可能导致震荡

## 更经济的版本：小批量随机梯度下降

原因：在整个训练集上计算开销过大。

随机采样 $b$ 个样本来近似损失。

$b$ :超参数，批量大小

$$\frac{1}{b} \sum_{i \in I_b} \ell(\mathbf{x}_i, y_i, \mathbf{w}) \quad (26)$$

批量不能过小，否则不能最大利用并行资源；

不能过大，增大开销，浪费计算。

## 动手做！

[linear](#)