

并行化计算尝试

2021 年 1 月 19 日

1 引言

给定一个矩阵 $X_{n \times m}$, 设 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ 是它的行向量, 我们想计算出这样一个矩阵 M , 满足 $M(i, j) = \cos(\mathbf{x}_i, \mathbf{x}_j)$, 具体有哪些计算方法呢?

2 方法一: 循环

首先是计算 M 的第一行, 然后是第二行, …… , 这样一行一行地计算:

```
1 import numpy as np
2 from time import time
3
4 # 要生成 1000 个随机点, 每行对应一个点, 每个点是一个 4 维向量.
5 n = 1000
6 data = np.random.rand(n, 4)
7
8 s = time()
9 cosines_method_1 = np.zeros((n, n))
10 for i in range(0, n-1):
11     for j in range(i+1, n):
12         inner_prod = data[i:(i+1), :] @ (data[j:(j+1), :].T)
13         norm_prod = np.linalg.norm(data[i, :]) * np.linalg.norm(data[j, :])
14         cosines_method_1[i, j] = inner_prod / norm_prod
15
16 np.fill_diagonal(cosines_method_1, 1)
17 e = time()
18 t_delta_1 = e-s
19 print(t_delta_1)
```

输出为

```
11.335773944854736
```

也就是说花了大概 11 秒的时间来计算这些.

3 方法二: 向量化

我们注意到, 对于每一组可能的 (i, j) , $M(i, j)$ 的计算都是相似的, 所以我们意识到, 有大量的相似的计算过程存在. 要将同样的计算方法应用在不同的元素上, 我们首先想到要利用线性代数教给我们的知识.

设 A 是一个 $n \times 1$ 的矩阵, 设 B 是一个 $1 \times n$ 的矩阵, 那么按照线性代数的知识, 我们知道 A 与 B 可以进行矩阵乘法操作, 也就是说, 存在矩阵 C , 它的行数等于 A 的行数, 它的列数等于 B 的列数, 并且对任意 $1 \leq i < j \leq n$, 都有 $C(i, j) = (AB)(i, j)$. 由此我们意识到, 如果一个列向量作为矩阵, 与自身的转置做矩阵乘法运算, 那么产生的那个矩阵的第 i 行第 j 列的元素的值, 就恰好等于它的第 i 个元素与第 j 个元素的乘积.

余弦值可以用范数与内积来表示, 设 \mathbf{a} 与 \mathbf{b} 是相同维数的向量, 那么

$$\cos \langle \mathbf{a}, \mathbf{b} \rangle = \frac{\mathbf{a} \mathbf{b}^T}{\|\mathbf{a}\| \|\mathbf{b}\|} \quad (3.1)$$

同样的, 对于 $M(i, j)$, 我们知道他是 $\cos \langle \mathbf{x}_i, \mathbf{x}_j \rangle$, 于是有

$$M(i, j) = \frac{\mathbf{x}_i \mathbf{x}_j^T}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \quad (3.2)$$

我们设想, 如果存在这样一个矩阵 A , 它满足: 对任意 $1 \leq i < j \leq n$, 都有 $A(i, j) = \mathbf{x}_i \mathbf{x}_j^T$; 并且如果还存在这样一个矩阵 B , 它满足, 对任意 $1 \leq i < j \leq n$, 都有 $B(i, j) = \|\mathbf{x}_i\| \|\mathbf{x}_j\|$, 那么矩阵 M 实际上就可以表示为 A 与 B 的按元素相除, 也就是说

$$M = A ./ B \quad (3.3)$$

现在, 我们应该开始考虑具体怎么计算 A 以及 B . 首先, 对于 A , 我们知道

$$A(i, j) = \mathbf{x}_i \mathbf{x}_j^T = \sum_{k=1}^m X(i, k) X(j, k) \quad (3.4)$$

为了简化问题, 我们假设 m 是一个不大的数, 比如说 $m = 4$, 这时

$$A(i, j) = X(i, 1)X(j, 1) + X(i, 2)X(j, 2) + X(i, 3)X(j, 3) + X(i, 4)X(j, 4) \quad (3.5)$$

我们故技重施, 将 A 写成四个矩阵的和

$$A = A_1 + A_2 + A_3 + A_4 \quad (3.6)$$

其中 $A_k(i, j) = X(i, k)X(j, k)$, 对于 $k = 1, 2, 3, 4; 1 \leq i < j \leq n$. 现在我们已经知道了, A_k 就是 X 的第 k 列与第 k 列转置的矩阵乘积, 也就是说

$$A_k = X(:, k) (X(:, k))^T, \quad k = 1, 2, 3, 4 \quad (3.7)$$

到了这里, 我们已经知道怎样把 A 的计算向量化了.

下面来看 B , 它的第 i 行第 j 列的元素的值正是 $\|\mathbf{x}_i\| \|\mathbf{x}_j\|$, 由此受到启发, 我们设想, 存在一个 $n \times 1$ 的矩阵 N , 满足 $N(i, 1) = \|\mathbf{x}_i\|$, 如果能计算出这个 N , 那么 M 就可以写为

$$M = N N^T \quad (3.8)$$

这个 N 也很好计算, 它等于

$$N = X(:, 1) \cdot X(:, 2) \cdot X(:, 3) \cdot X(:, 4) \quad (3.9)$$

也就是将 X 的四列看做四个矩阵, 再按元素相乘.

经过这一番分析, 我们已经将 M 拆解为一系列的矩阵运算与按元素运算的复合, 下面就可以开始实现了:

```

1 s= time()
2
3 data_1 = data[:, 0:1]
4 data_2 = data[:, 1:2]
5 data_3 = data[:, 2:3]
6 data_4 = data[:, 3:4]
7
8 pair_inner_prods = (data_1 @ data_1.T) + \
9     (data_2 @ data_2.T) + \
10    (data_3 @ data_3.T) + \
11    (data_4 @ data_4.T)
12
13 norms = np.sqrt(
14     data_1 * data_1 + \
15     data_2 * data_2 + \
16     data_3 * data_3 + \
17     data_4 * data_4
18 )
19
20 cosines_method_2 = pair_inner_prods / (norms @ norms.T)
21
22 e = time()
23 t_delta_2 = e-s
24 print(t_delta_2)

```

输出结果为

```
0.06543731689453125
```

再运行

```
abs(t_delta_1 - t_delta_2)/min(t_delta_1, t_delta_2)
```

可以看到快了约 170 倍。

4 方法三：GPU 加速

使用 NumPy 的矩阵乘法和按元素乘法功能之所以能大幅加快计算速度，是因为 NumPy 利用了专业的线性代数数值计算库如 BLAS, LAPACK 等，而这些数值计算库则充分利用了 CPU 的 SIMD¹ 指令，简而言之，NumPy 让我们在这个计算任务上，充分发挥了 CPU 的大部分潜力。

而对于 SIMD 类型的计算任务，GPU 比 CPU 其实是更加擅长的，再加上现在有了 CUDA² 和 CuPy³，想要验证这一点变得非常简单。

首先载入 CuPy 软件包，并且将数据迁移到显存上：

```

1 import cupy as cp
2
3 s_move_data_to_gpu = time()
4 data_1 = cp.asarray(data_1)
5 data_2 = cp.asarray(data_2)

```

¹Single Instruction Multiple Data, 是指一系列可同时操作于非常多个运算数的运算指令

²Nvidia 公司为它家特定系列的显卡开发的一系列并行计算 API

³CUDA 提供给 Python 开发者的编程接口

```

6 data_3 = cp.asarray(data_3)
7 data_4 = cp.asarray(data_4)
8 e_move_data_to_gpu = time()
9 print(e_move_data_to_gpu - s_move_data_to_gpu)

```

然后就可以开始在 GPU 上进行计算了:

```

1 s_gpu_compute = time()
2
3 pair_inner_prods = cp.matmul(data_1, data_1.T) + \
4     cp.matmul(data_2, data_2.T) + \
5     cp.matmul(data_3, data_3.T) + \
6     cp.matmul(data_4, data_4.T)
7
8 norms = cp.sqrt(
9     data_1 * data_1 + \
10    data_2 * data_2 + \
11    data_3 * data_3 + \
12    data_4 * data_4
13 )
14
15 cosines_method_2 = pair_inner_prods / cp.matmul(norms, norms.T)
16
17 e_gpu_compute = time()
18
19 print(e_gpu_compute - s_gpu_compute)

```

输出结果为:

```
0.0031576156616210938
```

事实上 CuPy 与 NumPy 是很大程度上兼容的, 运行了 `cp.asarray` 之后, 剩下的无非就是把代码里的 `np` 替换为 `cp`, 就可以通过调用 CuPy 在 GPU 上进行计算了. 现在再运行:

```
t_delta_2 / (e_gpu_compute - s_gpu_compute)
```

可看到输出

```
20.723648444612217
```

这说明 GPU 比 CPU 还要再快上一个台阶.

5 结论

我们先后尝试了三种计算余弦矩阵 M 的方法, 第一种方法是用 `for` 循环自己编程实现, 耗时为 11.335773, 第二种方法是把原有的计算任务用向量与矩阵之间的操作来实现, 耗时为 0.065437, 第三种方法是把第二种方法放到 GPU 上进行, 耗时为 0.003157. 通过简单计算可得, 第二种方法约比第一种快了 170 倍, 第三种方法又比第一种方法快了 20 倍, 可想而知, 第三种方法比第一种方法快了数千倍.