

Cấu Trúc Dữ Liệu Và Giải Thuật

Bài Tập Lớn 1

XỬ LÝ DỮ LIỆU MNIST VÀ HIỆN THỰC GIẢI THUẬT kNN BẰNG CẤU TRÚC DỮ LIỆU DANH SÁCH

nhóm thảo luận CSE
<https://www.facebook.com/groups/211867931379013>

Tp. Hồ Chí Minh, Tháng 2/2024



Mục lục

1	Kế hoạch làm bài	3
2	Task 2: Hiện thực Dataset	4
2.1	Lý Thuyết	4
2.1.1	Phần dữ liệu X (gồm X_{train} và X_{test})	5
2.1.2	Phần dữ liệu Y (gồm Y_{train} và Y_{test})	5
2.2	Một số Hàm dùng trong BTL này	6
2.2.1	Bài Tập Lớn	8
3	Các Khóa Học HK232	13



1 Kế hoạch làm bài

Phần Làm	Bắt đầu	Thời gian	Kết thúc	Kiến thức	Nộp lại discord anh
Task1	26/2	6 ngày	6:00 2/3	Danh sách liên kết	ID_task1
Task2	2/3	6 ngày	6:00 7/3	Danh sách liên kết	ID_task2
Task3(full)	7/3	6 ngày	6:00 13/3	Danh sách liên kết	ID_task3

Một số quy định của nhóm

- Nếu mấy bạn vô trễ thì thời gian làm có thể kéo dài theo ngày với thời gian Bắt đầu sẽ là ngày hôm đó
- Sau *deadline* thì anh không chỉ BTL phần đó nữa 😊 😊 😊
- Bạn nào trễ hạn 1 lần BTL2 anh sẽ không giảm giá
- Hỏi lý thuyết thì hỏi trong nhóm cho tiện
- Hỏi code thì nhắn riêng anh TA để bảo mật cho mấy bạn có thông báo trong nhóm
- Nộp bài thì gửi file code *kNN.cpp* và *kNN.hpp* chuyển thành tên *ID_task2.cpp* và *ID_task2.hpp* với *ID* trong excel nộp qua anh TA luôn



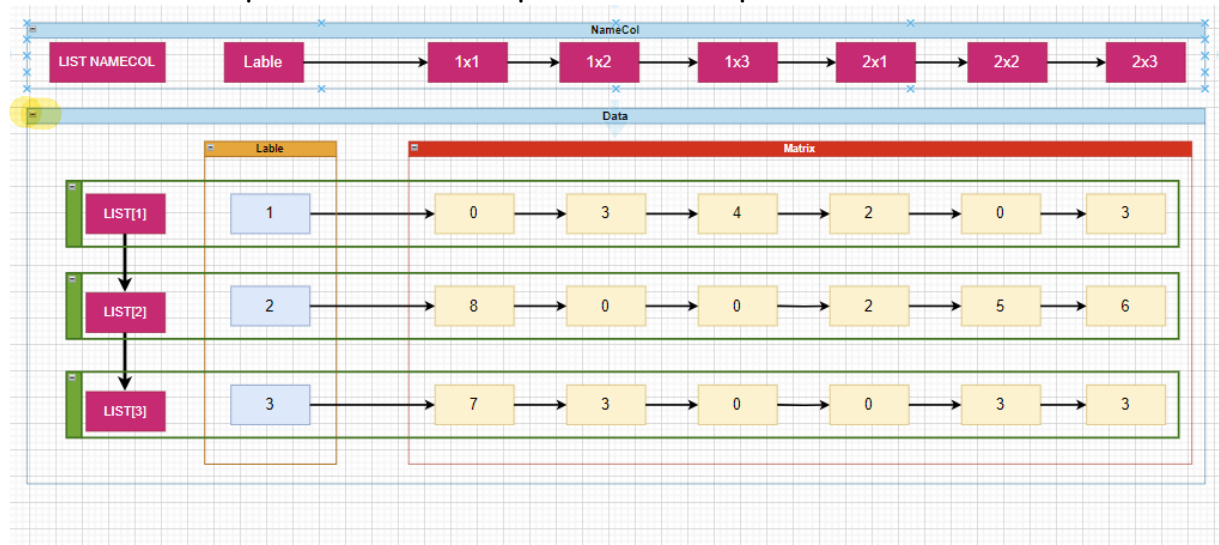
2 Task 2: Hiện thực Dataset

2.1 Lý Thuyết

Sinh viên được cung cấp file `mnist.csv`, trong đó:

- Mỗi dòng là mô tả một hình ảnh chữ số, được mã hóa thành $28 \times 28 = 784$ cột và một cột nhãn.
- Cột nhãn (label) đánh nhãn cho số được thể hiện bởi hình ảnh. Đối với mỗi ảnh, khoảng giá trị của ô trong cột này là từ 0 đến 9.
- Các cột còn lại, có tiêu đề `ixj` trong đó i, j nằm trong khoảng $[1, 28]$ thể hiện tọa độ của điểm ảnh. Đối với mỗi ảnh, khoảng giá trị của ô trong cột này là từ 0 đến 255 (độ đậm của điểm ảnh).

Sau khi đưa dữ liệu vào *Dataset* ta được cấu trúc dữ liệu sau



Hình 1: Dataset ban đầu

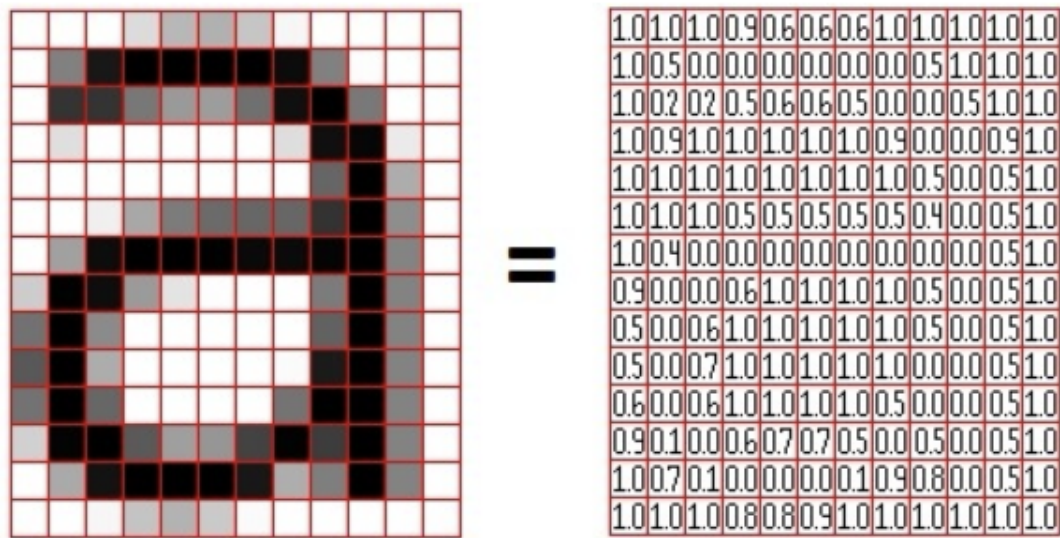
- NameCol** : là hàng đầu tiên trong `mnist.csv` là tên của của cột trong bảng được thiết kế như một danh sách liên kết được dùng ở bước trước

```

1 List<string>* nameCol;
2 // khởi tạo
3 nameCol = new Image<string>();
4 // các hàm hay dùng trong bài này đối với nameCol
5 List<T>* subList(int start, int end)
6 void printStartToEnd(int start, int end) const
7 void print() const
8 void push_back(T value)
9 int length() const
10 T& get(int index) const

```

- Data** : là phần còn lại của trong `mnist.csv` là danh sách ma trận 2 chiều với chiều được thiết kế bằng danh sách liên kết như hình vẽ với mỗi phần `List[N]` trong `Data` là một danh sách liên kết con
 - `Lable` là kết quả của hình ảnh đang là số 1 đến số 9 phần sau ta sẽ cắt nó ra thành phần `Y`
 - `Matrix` là dữ liệu của ảnh đã được mã hóa thành một ma trận hai chiều đã được mã hóa thành một ma trận 1 chiều (Xem cách chuyển array 2 chiều thành array) phần sau ta sẽ cắt nó ra thành phần `X`



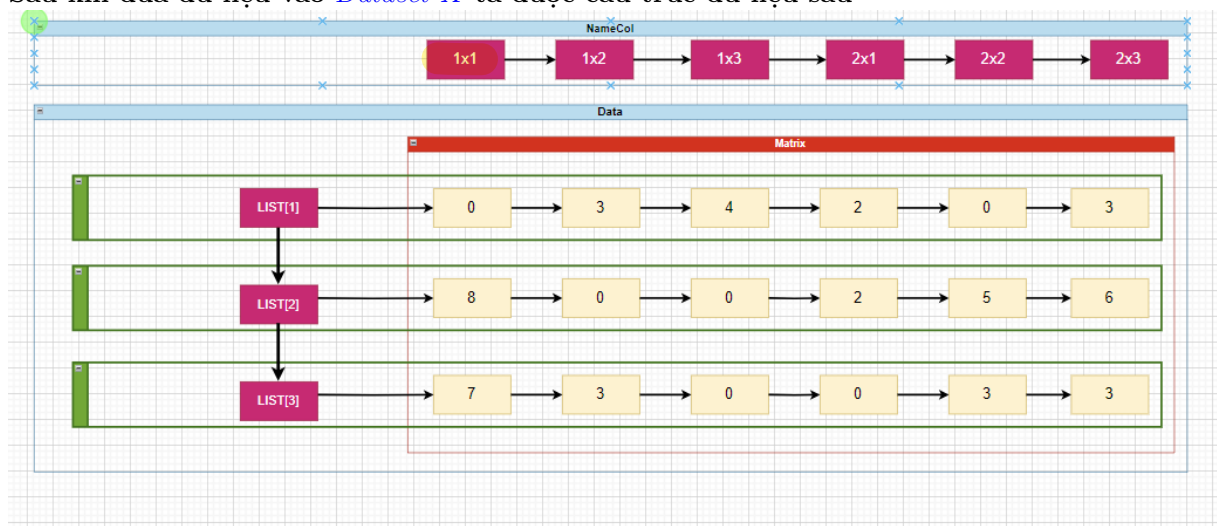
Hình 2: Hình ảnh khi chuyển về ma trận với $lable = a$ và $matrix$ như hình bên phải

2.1.1 Phần dữ liệu X (gồm X_{train} và X_{test})

Phần dữ liệu X là phần $matrix$ phần này ta chưa xác định được ảnh là số nằm trong khoảng $[1,9]$

- X_{train} Bảng dữ liệu huấn luyện chứa các đặc trưng của ảnh. Số hàng là số lượng ảnh huấn luyện, Số cột là số lượng các đặc trưng. có thể hiểu là phần này là phần dùng để cho con AI nó học thông minh lên
- X_{test} Bảng dữ liệu dự đoán chứa các đặc trưng của ảnh. Số hàng là số lượng ảnh, Số cột là số lượng các đặc trưng. có thể hiểu là phần này dùng để xem con AI nó có ngu hay không bằng cách đưa ảnh vào xem nó đoán được đúng số hay không

Sau khi đưa dữ liệu vào $Dataset X$ ta được cấu trúc dữ liệu sau



Hình 3: $Dataset X$

2.1.2 Phần dữ liệu Y (gồm Y_{train} và Y_{test})

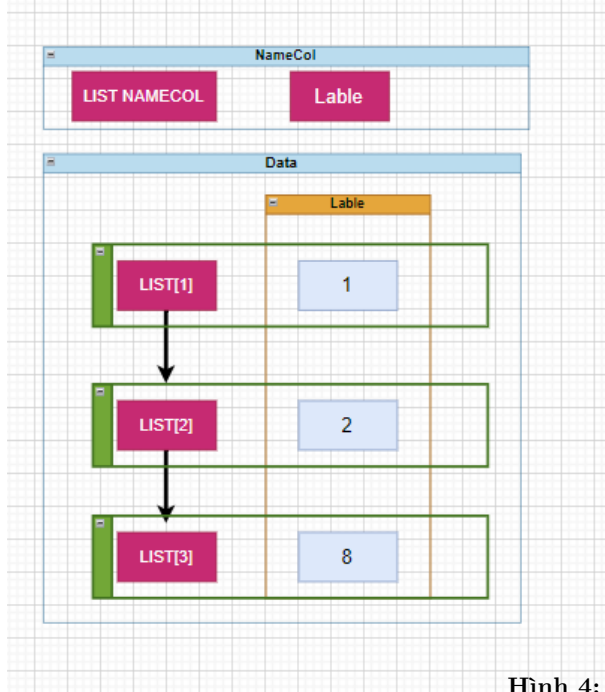
Phần dữ liệu Y là phần $lable$ phần này ta xác định được ảnh là số nằm trong khoảng $[1,9]$

- Y_{train} Bảng dữ liệu huấn luyện chứa nhãn của ảnh. Số hàng là số lượng ảnh huấn luyện, Số cột là 1. có thể hiểu là phần này là phần dùng để cho con AI nó học thông minh lên



- *Y_test* Bảng dữ liệu dự đoán chứa nhãn thực sự của ảnh. Số hàng là số lượng ảnh, Số cột là 1. **có thể hiểu là phần này dùng để xem con AI nó có ngu hay không** bằng cách đưa ảnh vào xem nó đoán được đúng số hay không

Sau khi đưa dữ liệu vào *Dataset Y* ta được cấu trúc dữ liệu sau



Hình 4: *Dataset Y*

2.2 Một số Hàm dùng trong BTL này

1. Khởi tạo và hủy *data*

```
1 // khởi tạo
2 this->nameCol = new Image<string>();
3 // hủy
4 delete data;
```

2. Lấy chiều rộng và chiều dài Chú ý trường hợp rỗng

```
1 nRows = this->data->length();
2 nCols = this->data->get(0)->length(); // không xảy ra với Th nRow = 0
3
4 // kết quả với dataset X
5 nRows = 3
6 nCols = 6
7
8 // kết quả với dataset Y
9 nRows = 1 // phần này đối với Y là cố định
10 nCols = 1
```

3. Print ra trong khoảng từ *[start, end]*



```
1 // print ra 5 phần tử hàng đầu tiên
2 data->get(0)->printStartToEnd(0, 5);
3 // kết quả với dataset X
4 0 3 7 5 0
5 // kết quả với dataset Y
6 1
7
8 // print ra 1 phần tử hàng thứ 3
9 data->get(2)->printStartToEnd(0, 1);
10 // kết quả với dataset X
11 7
12 // kết quả với dataset Y
13 8
14
15 // print ra tất cả hàng thứ 2
16 data->get(1)->print();
17 // kết quả với dataset X
18 8 0 0 2 5 6
19 // kết quả với dataset Y
20 2
```

4. Truy cập đến phần bất kì trong *matrix*

```
1 // lấy ra hàng thứ 2 là danh sách liên kết
2 data->get(1)
3
4 // lấy ra phần tử thứ 3 trong hàng thứ 3
5 data->get(2)->get(2)
```

5. Xóa một hàng hoặc một cột bất kì trong một *matrix*

```
1 // xóa đi hàng đầu tiên
2 data->remove(0)
3
4 // xóa đi cột thứ 2
5 for(int i = 0; i < row; i++) data->get(i)->remove(1);
6
7 // xóa phần tử hàng thứ 3 trong cột thứ 3
8 data->get(2)->remove(2);
```

6. Cắt ma trận

```
1 // lấy 2 hàng đầu tiên và 2 cột đầu tiên trong ma trận
2 Dataset result;
3 result.nameCol = nameCol->subList(0, 2);
4 for(int i = 0; i < 2; i++)
5 {
6     result.data->push_back(data->get(i)->subList(0, 2));
7 }
```



2.2.1 Bài Tập Lớn

```
1 class Dataset {
2 private:
3     List<List<int>*>* data;
4     List<string>* NameCol;
5     //You may need to define more
6 public:
7     //! Hàm khởi tạo
8     Dataset::Dataset()
9     {
10         this->NameCol = new Image<string>();
11         this->data = new Image<List<int>*>();
12     }
13     //! Hàm hủy
14     Dataset::~Dataset()
15     {
16         delete data;
17         delete nameCol;
18     }
19     Dataset::Dataset(const Dataset& other);
20     Dataset& Dataset::operator=(const Dataset& other);
21     bool Dataset::loadFromCSV(const char* fileName);
22     void Dataset::getShape(int& nRows, int& nCols) const;
23     void Dataset::columns() const;
24     void printHead(int nRows = 5, int nCols = 5) const;
25     void printTail(int nRows = 5, int nCols = 5) const;
26     bool drop(int axis = 0, int index = 0, std::string columns = "");
27     Dataset extract(int startRow = 0, int endRow = -1, int startCol = 0, int endCol
28         → = -1) const;
29
30     Dataset predict(const Dataset& X_train, const Dataset& Y_train, const int k)
31         → const;
32     double score(const Dataset& y_test) const;
33 };
```

1. **data** và **nameCol** đã được định nghĩa trong phần lý thuyết
2. Hàm khởi tạo rỗng sẽ gán các giá trị **data** và **nameCol** cấp phát bộ nhớ trong vùng nhớ **heap**
3. Hàm hủy là hủy đi **data** và **nameCol**
4. Hàm khởi tạo có tham số đầu vào **other** hàm này yêu cầu **copy** đối tượng **other** sang cho đối tượng cần được khởi tạo
5. Hàm **operator=** hàm này dùng để **copy** giống hàm trên
6. Hàm **loadFromCSV** khi chạy hết test case các bạn nên hiện thực lại không bị chấm đạo văn, Đoạn code trong đó chỉ mang tính tham khảo



7. Hàm **getShape(int& nRows, int& nCols)** gán *nRows* và *nCols* là chiều dài và chiều rộng của *matrix*

```
1 X_train.getShape(nRows, nCols)
2 // kết quả
3 nRows = 3
4 nCols = 6
5
6 Y_train.getShape(nRows, nCols)
7 // kết quả
8 nRows = 3
9 nCols = 1 // luôn luôn = 1 chỉ có một cột lable
```

8. Hàm **columns()** in ra tất cả các tên

```
1 X_train.columns();
2 // kết quả
3 1x1 1x2 1x3 2x1 2x2 2x3
4
5 Y_train.columns();
6 // kết quả
7 lable
```

9. Hàm **printHead(int nRows = 5, int nCols = 5)** In ra nRows dòng đầu tiên, và chỉ in nCols cột đầu tiên của bảng dữ liệu.

- Dòng đầu tiên, in ra tên các cột của bảng dữ liệu (dùng *nameCol*)
- Từ dòng thứ 2 trở đi, in giá trị của mỗi ô trong bảng, mỗi phần tử cách nhau bằng khoảng trắng, không có khoảng trắng dư ở cuối dòng (dùng *data* và hàm *printStartToEnd* để xử lý)
- Nếu *nRows* lớn hơn số lượng dòng trong bảng dữ liệu, in hết dòng trong bảng dữ liệu. Nếu *nCols* lớn hơn số lượng cột trong bảng dữ liệu, in hết cột trong bảng dữ liệu. Nếu *nRows* hoặc *nCols* nhỏ hơn 0, không in gì cả.

```
1 X_train.printHead(2, -1);
2 // kết quả
3
4 X_train.printHead(2, 2);
5 // kết quả
6 1x1 2x2
7 0 3
8 8 0
9
10 Y_train.printHead(2, 2);
11 // kết quả
12 lable
13 1
14 2
```

10. Hàm **printTail(int nRows = 5, int nCols = 5)** In ra nRows dòng cuối cùng, và chỉ in nCols cột cuối cùng của bảng dữ liệu..

- Dòng đầu tiên, in ra tên các cột của bảng dữ liệu (dùng *nameCol*)



- Từ dòng thứ 2 trở đi, in giá trị của mỗi ô trong bảng, mỗi phần tử cách nhau bằng khoảng trắng, không có khoảng trắng dư ở cuối dòng (dùng *data* và hàm *printStartToEnd* để xử lý)
- Nếu nRows lớn hơn số lượng dòng trong bảng dữ liệu, in hết dòng trong bảng dữ liệu. Nếu nCols lớn hơn số lượng cột trong bảng dữ liệu, in hết cột trong bảng dữ liệu. Nếu nRows hoặc nCols nhỏ hơn 0, không in gì cả.

```

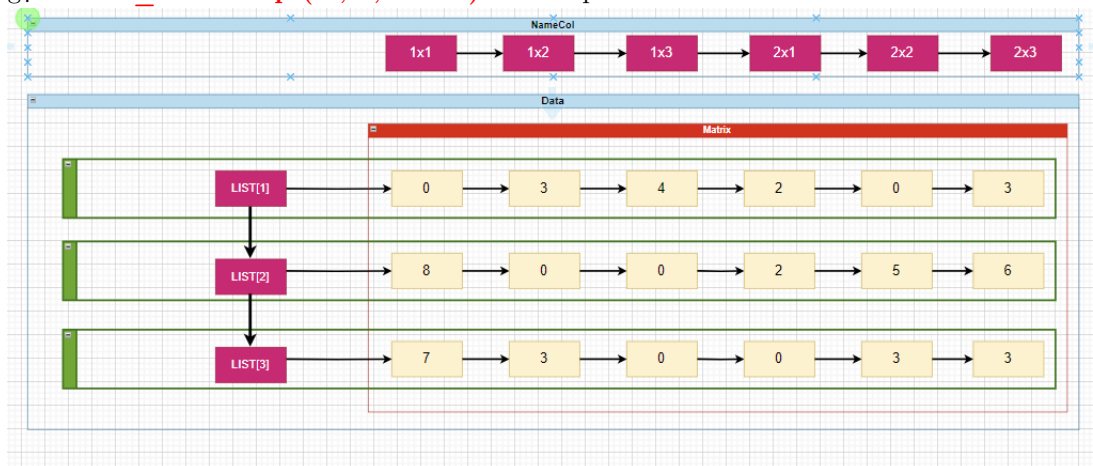
1 X_train.printTail(2, -1);
2 // kết quả
3
4 X_train.printTail(2, 2);
5 // kết quả
6 2x2 2x3
7 5 6
8 3 3
9
10 Y_train.printTail(2, 2);
11 // kết quả
12 lable
13 2
14 8

```

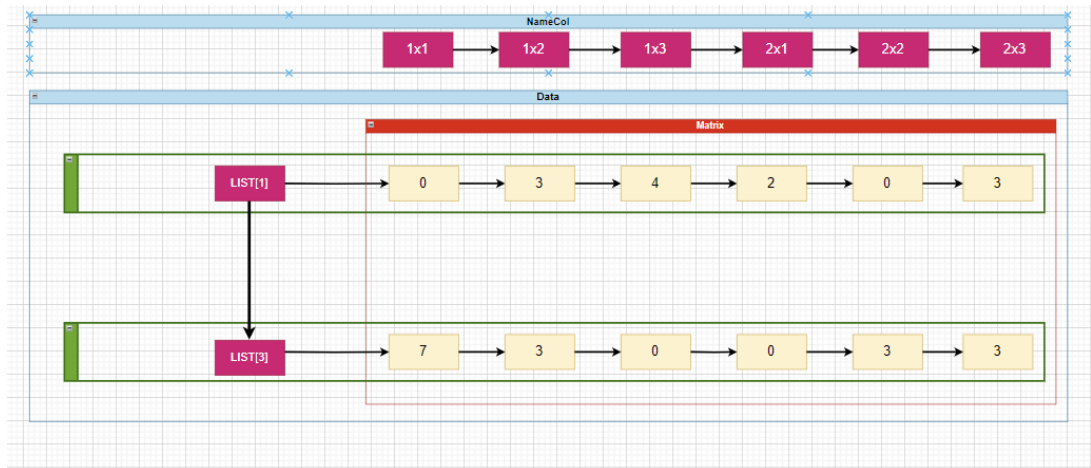
11. Hàm `drop (int axis = 0, int index = 0, std::string columns = "")`

- Xóa một hàng hoặc cột của bảng dữ liệu.
- Nếu axis khác 0 hoặc 1, hàm không làm gì cả và trả về false.
- Nếu axis == 0, thực hiện xóa một hàng tại vị trí index, sau đó trả về true. Nếu index >= số dòng hoặc < 0, hàm không làm gì cả và trả về false.
- Nếu axis == 1, thực hiện xóa một cột có tên trùng với columnName, sau đó trả về true. Nếu columnName không nằm trong danh sách tên các cột, hàm không làm gì cả và trả về false.

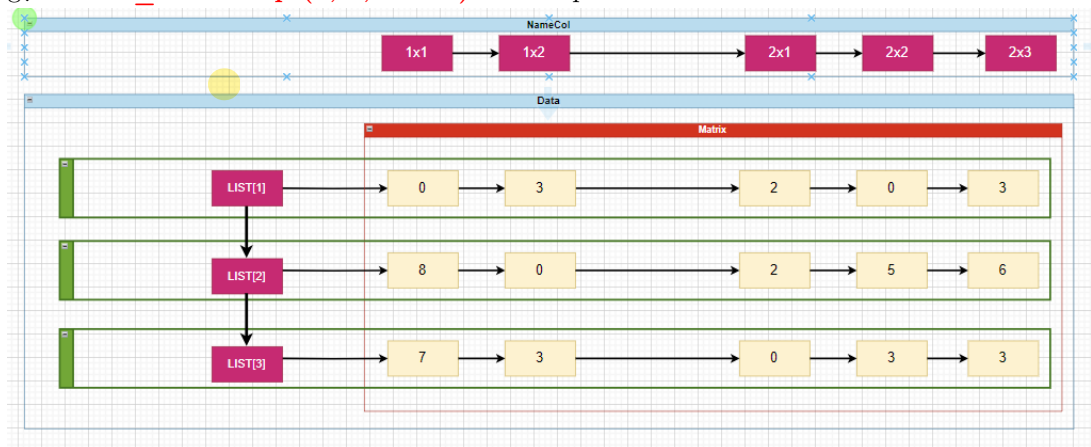
(a) gọi hàm `X_train.drop (-1, 0, "1x3")` thì kết quả sẽ như hình bên dưới



(b) gọi hàm `X_train.drop (0, 1, "1x3")` thì kết quả sẽ như hình bên dưới



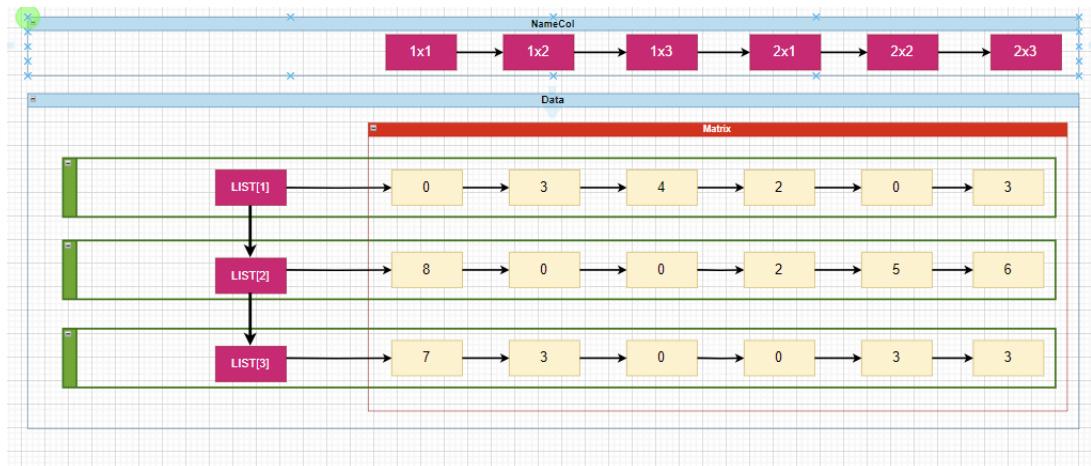
(c) gọi hàm `X_train.drop (1, 1, "1x3")` thì kết quả sẽ như hình bên dưới



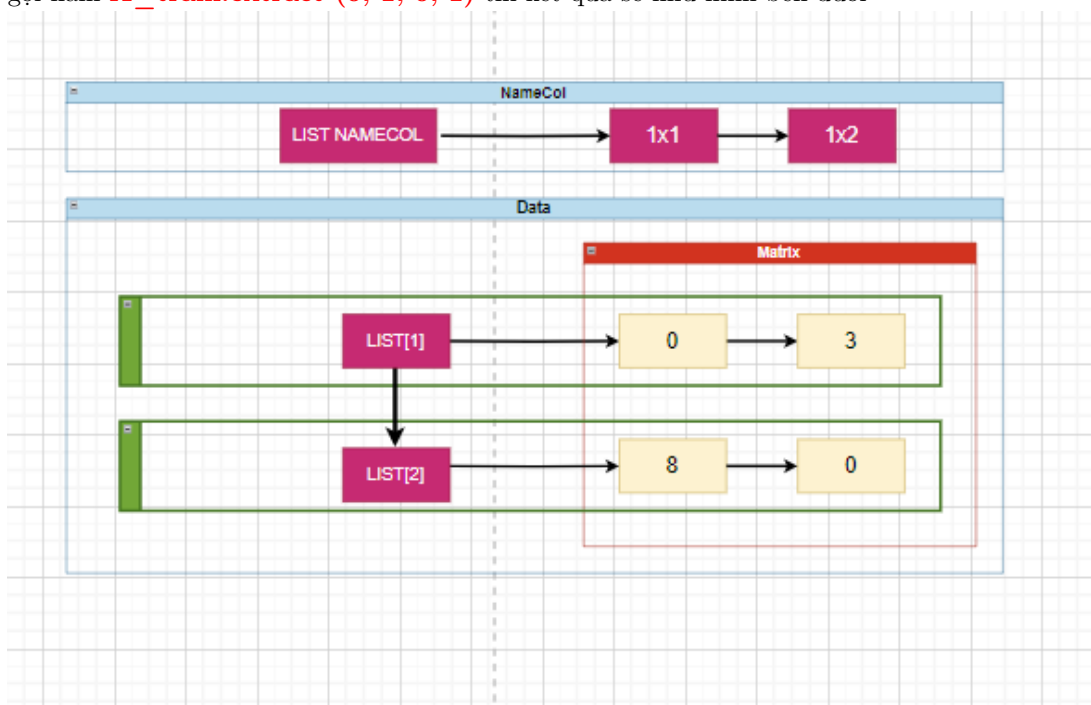
12. Hàm `extract(int startRow = 0, int endRow = -1, int startCol = 0, int endCol = -1)`

- Phương thức dùng để trích xuất một phần của bảng dữ liệu, sau đó trả về bảng dữ liệu đã trích xuất.
- Trong đó: `startRow` là hàng bắt đầu, `endRow` là hàng cuối cùng, `startCol` là cột bắt đầu, `endCol` là cột cuối cùng.
- Nếu `endRow == -1`, ta sẽ lấy tất cả các hàng. Tương tự với `endCol == -1`. và trong trường hợp `endRow`, `endCol` lớn hơn kích thước thì cũng lấy hết tới cuối
- Các testcase sẽ đảm bảo `startRow`, `endRow`, `startCol` và `endCol` nằm trong khoảng giá trị hợp lệ (từ 0 đến số hàng/số cột) và `start <= end`

(a) gọi hàm `X_train.extract (0,-1,0,-1)` thì kết quả sẽ như hình bên dưới



(b) gọi hàm **X_train.extract (0, 1, 0, 1)** thì kết quả sẽ như hình bên dưới



13. Hàm **distanceEuclidean(const List<int>* x, const List<int>* y)** Copy từ class *Image* xuống và chỉnh lại. Và bỏ *distanceEuclidean* trong class *Image*



3 Các Khóa Học HK232

nhóm thảo luận CSE

<https://www.facebook.com/groups/211867931379013>

- Lớp BTL1 + GK + LAB + Lý thuyết + Harmony của môn DSA HK232
- Lớp BTL2 + CK + LAB + Lý thuyết + Harmony của môn DSA HK232
- Lớp BTL1 + LAB + Lý thuyết + Harmony của môn KTLT HK232
- Lớp BTL2 + LAB + Lý thuyết + Harmony của môn KTLT HK232
- Lớp BTL1 + BTL2 + GK + Harmony của môn PPL HK232
- Lớp BTL3 + BTL4 + CK + Harmony của môn PPL HK232

CHÚC CÁC EM HỌC TỐT

