# Feature selection in high-dimensional dataset using MapReduce

**Claudio Reggiani**

Yann-Aël Le Borgne

Gianluca Bontempi

BNAIC 2017 – November 8, 2017

mRMR

# mRMR

$$\max_{x_j \in X - S_{m-1}}$$

candidate
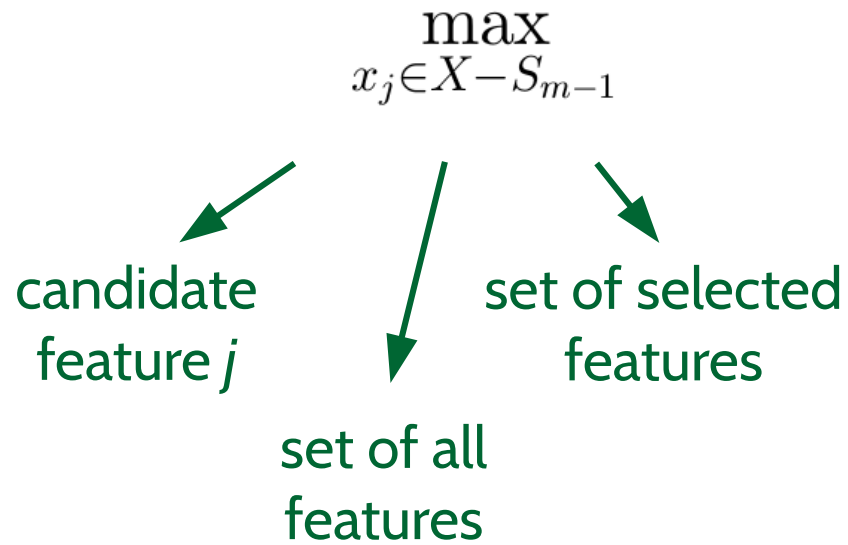feature $j$

# mRMR

$$\max_{x_j \in X - S_{m-1}}$$

candidate
feature *j*

set of all
features

# mRMR

$$\max_{x_j \in X - S_{m-1}}$$

candidate
feature $j$

set of all
features

set of selected
features

# mRMR

minimal redundancy

$$\max_{x_j \in X - S_{m-1}} \left[ \phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxx} \right]$$

candidate
feature $j$

set of all
features

set of selected
features

# mRMR

$$\max_{x_j \in X - S_{m-1}} \left[ \overbrace{I(x_j; c)}^{\text{maximal relevance}} \overbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxx}}^{\text{minimal redundancy}} \right]$$

candidate feature $j$

set of all features

set of selected features

# mRMR

$$\max_{x_j \in X - S_{m-1}} \left[ \overbrace{I(x_j; c)}^{\text{maximal relevance}} - \overbrace{\sum_{x_i \in S_{m-1}} I(x_j; x_i)}^{\text{minimal redundancy}} \right]$$

candidate feature *j*

set of all features

set of selected features

# mRMR



maximal relevance

minimal redundancy

$$\max_{x_j \in X - S_{m-1}} \left[ I(x_j; c) - \sum_{x_i \in S_{m-1}} I(x_j; x_i) \right]$$

candidate feature $j$

set of all features

set of selected features

mutual information
(descrete values)

# mRMR

maximal relevance · minimal redundancy

$$\max_{x_j \in X - S_{m-1}} \left[ I(x_j; c) - \frac{1}{m-1} \sum_{x_i \in S_{m-1}} I(x_j; x_i) \right]$$

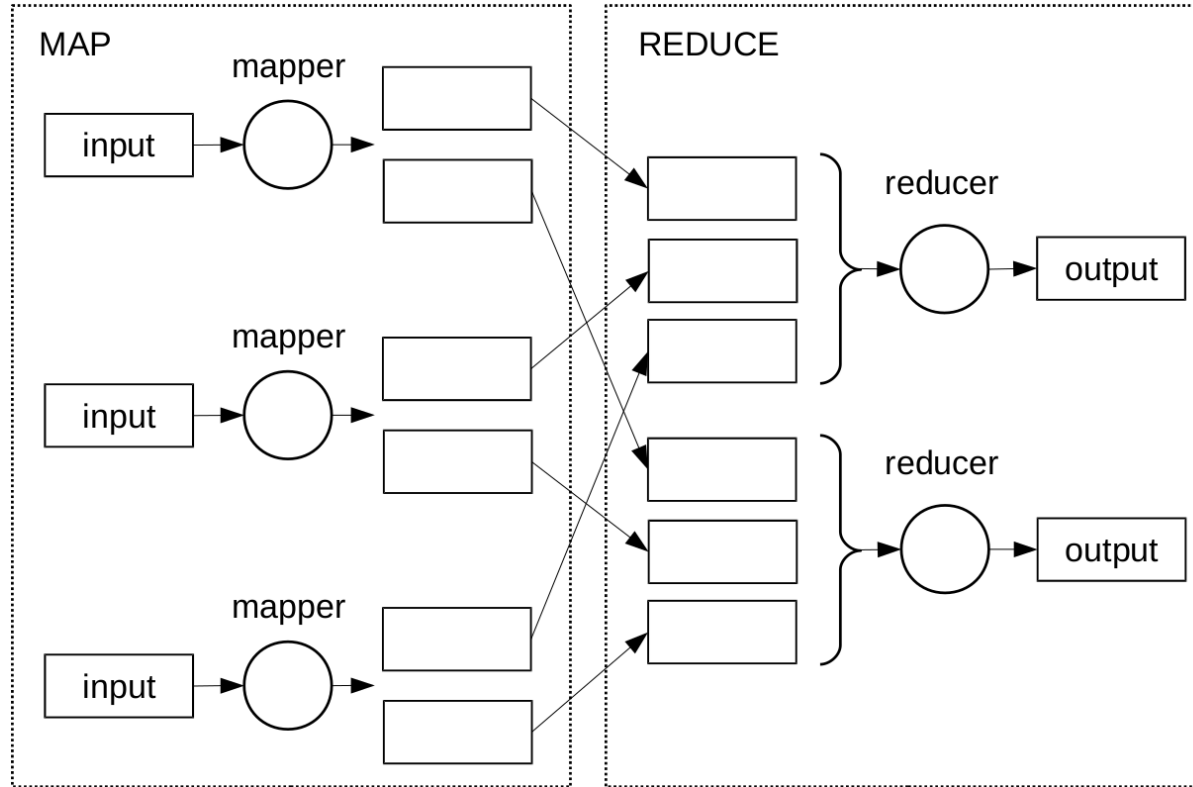candidate feature $j$

set of all features

set of selected features

mutual information (descrete values)

# mRMR

maximal relevance

minimal redundancy

$$\max_{x_j \in X - S_{m-1}} \left[ I(x_j; c) - \frac{1}{m-1} \sum_{x_i \in S_{m-1}} I(x_j; x_i) \right]$$

candidate feature $j$

set of all features

set of selected features

cardinality of selected features set

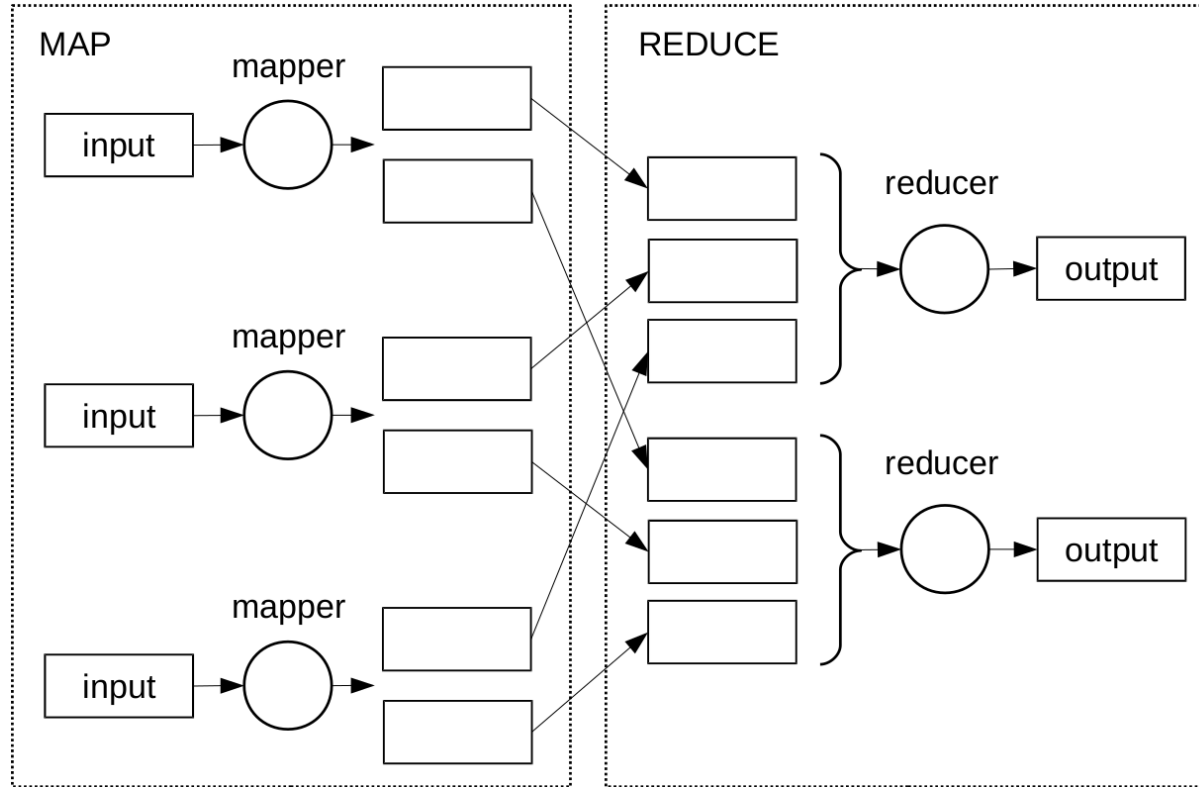mutual information (descrete values)

How do we cope with high-dimensional datasets?
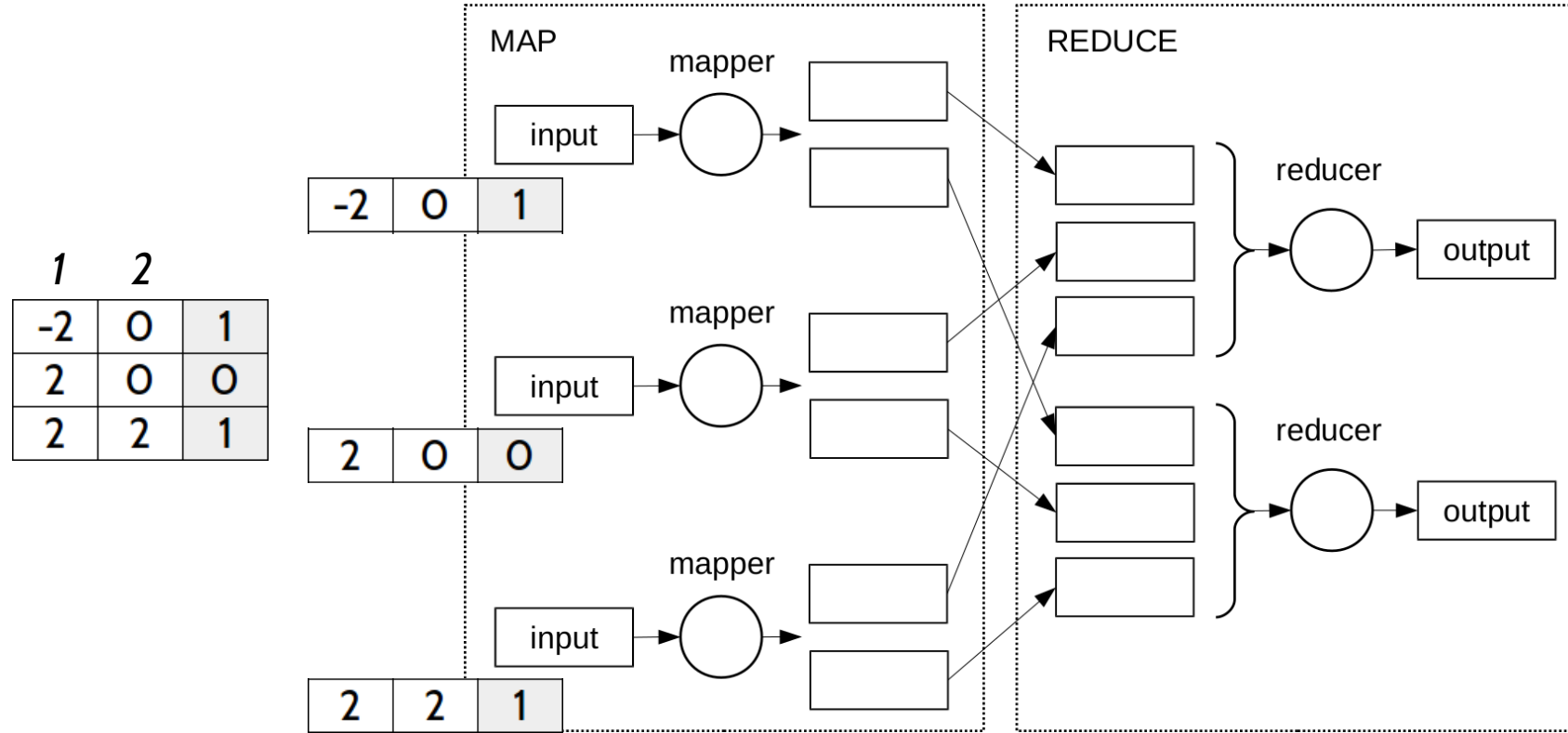
# Brief introduction to MapReduce
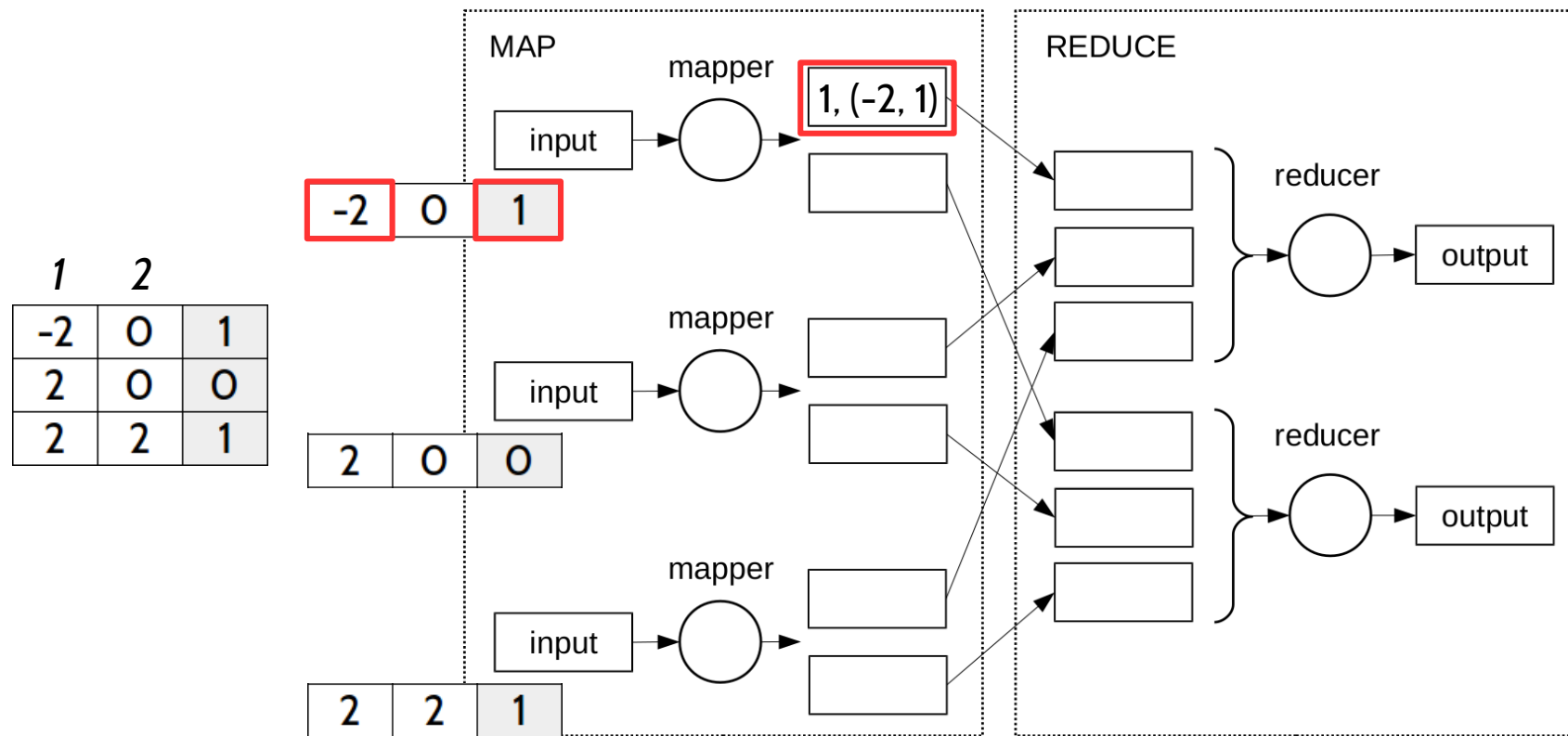
# Brief introduction to MapReduce

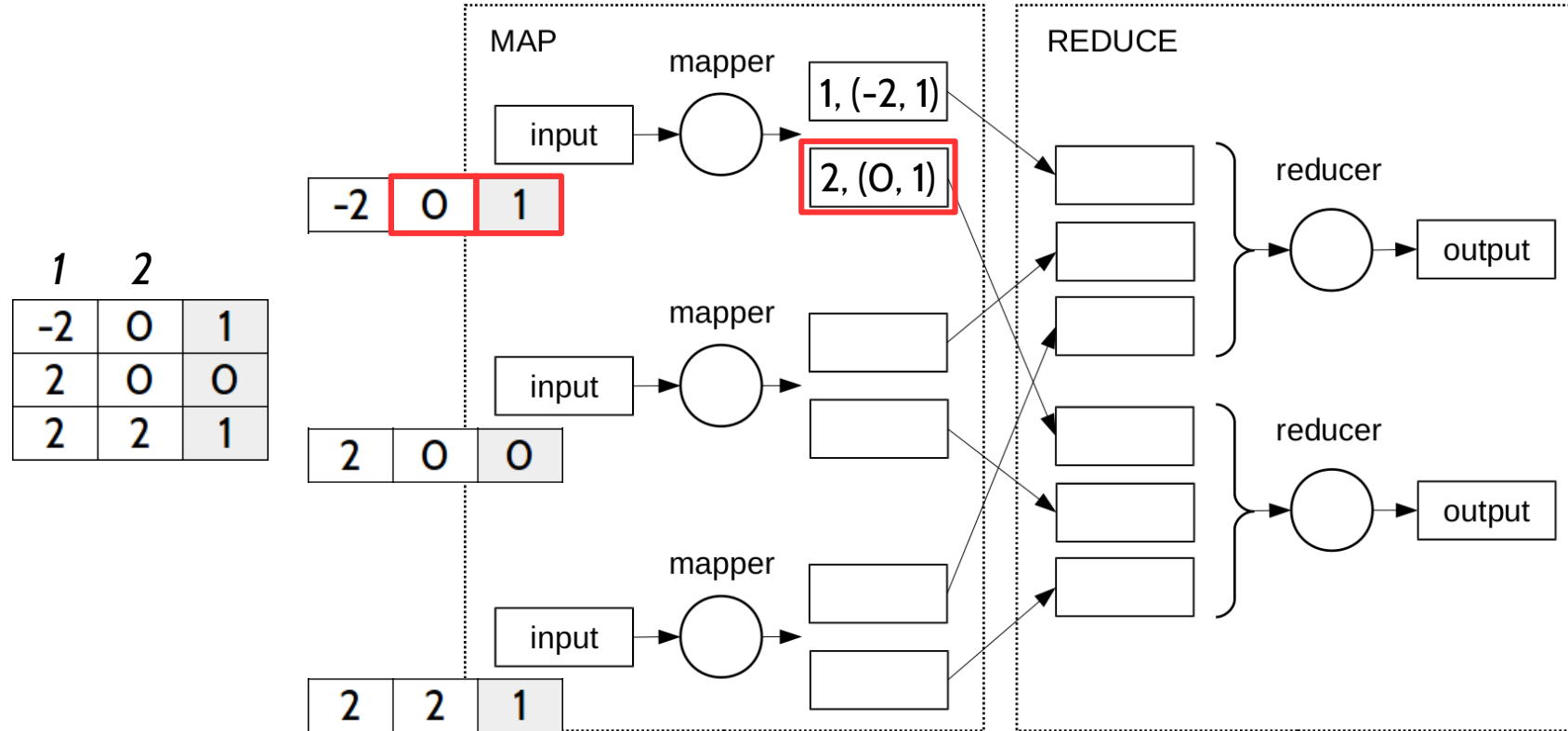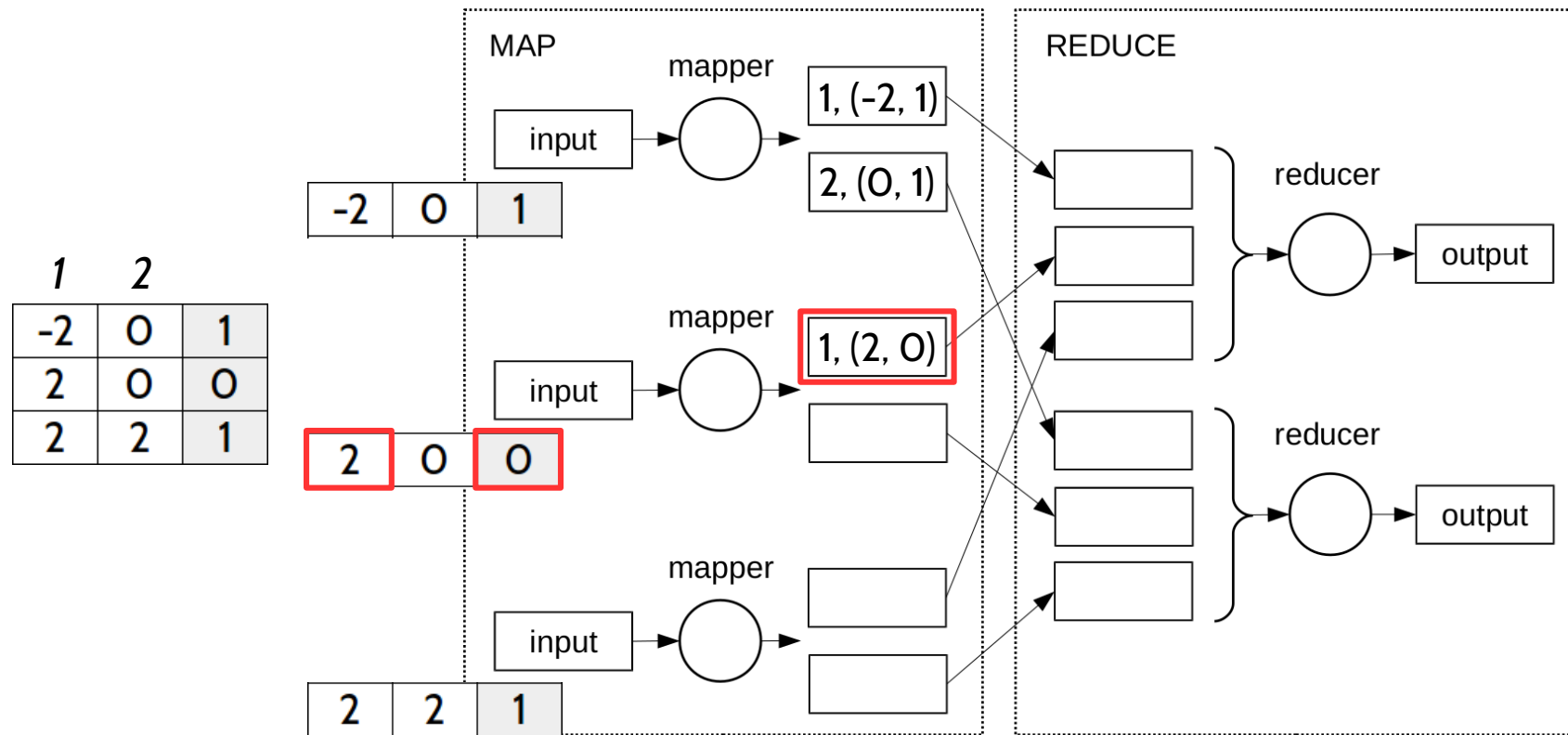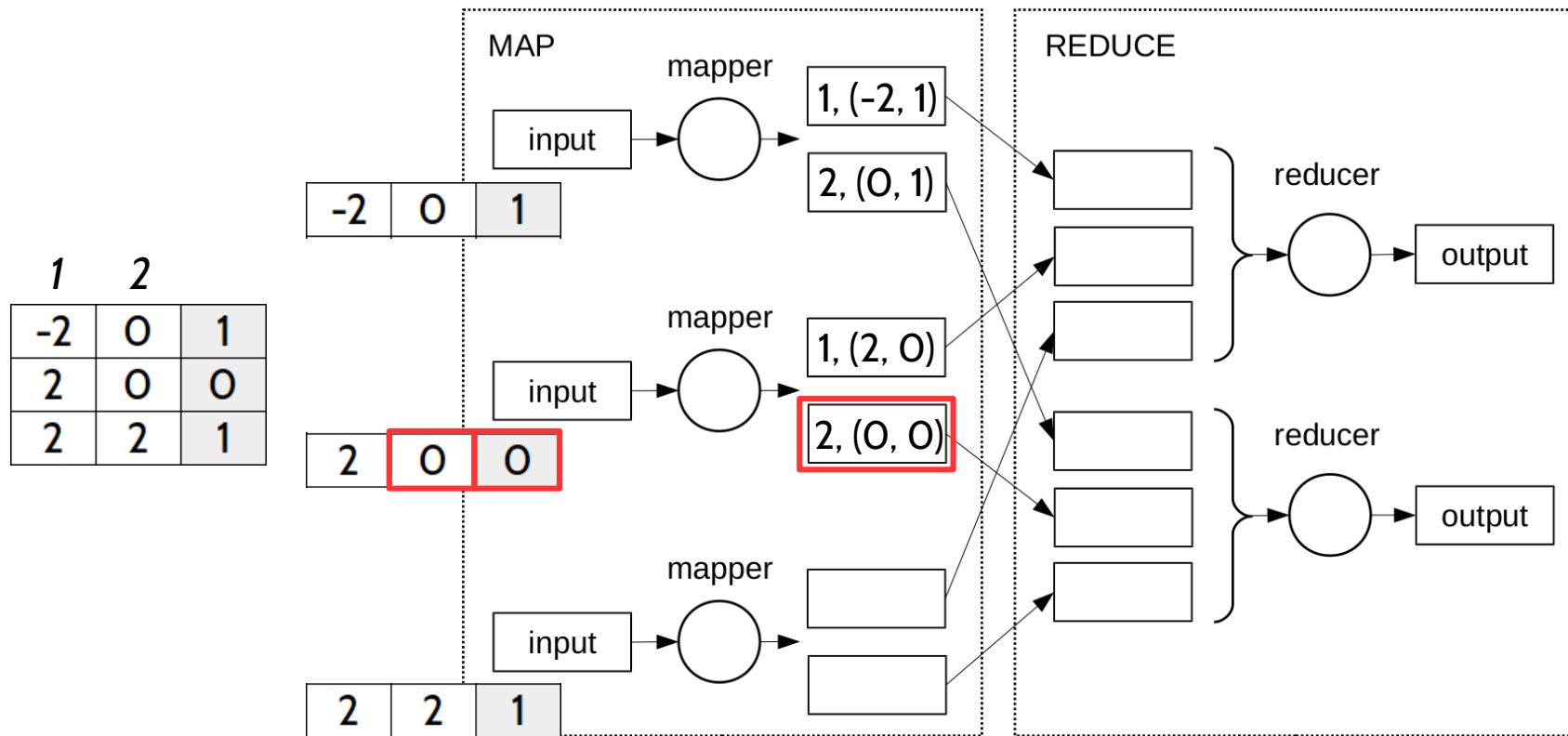| *1* | *2* | |
|---|---|---|
| -2 | 0 | 1 |
| 2 | 0 | 0 |
| 2 | 2 | 1 |

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

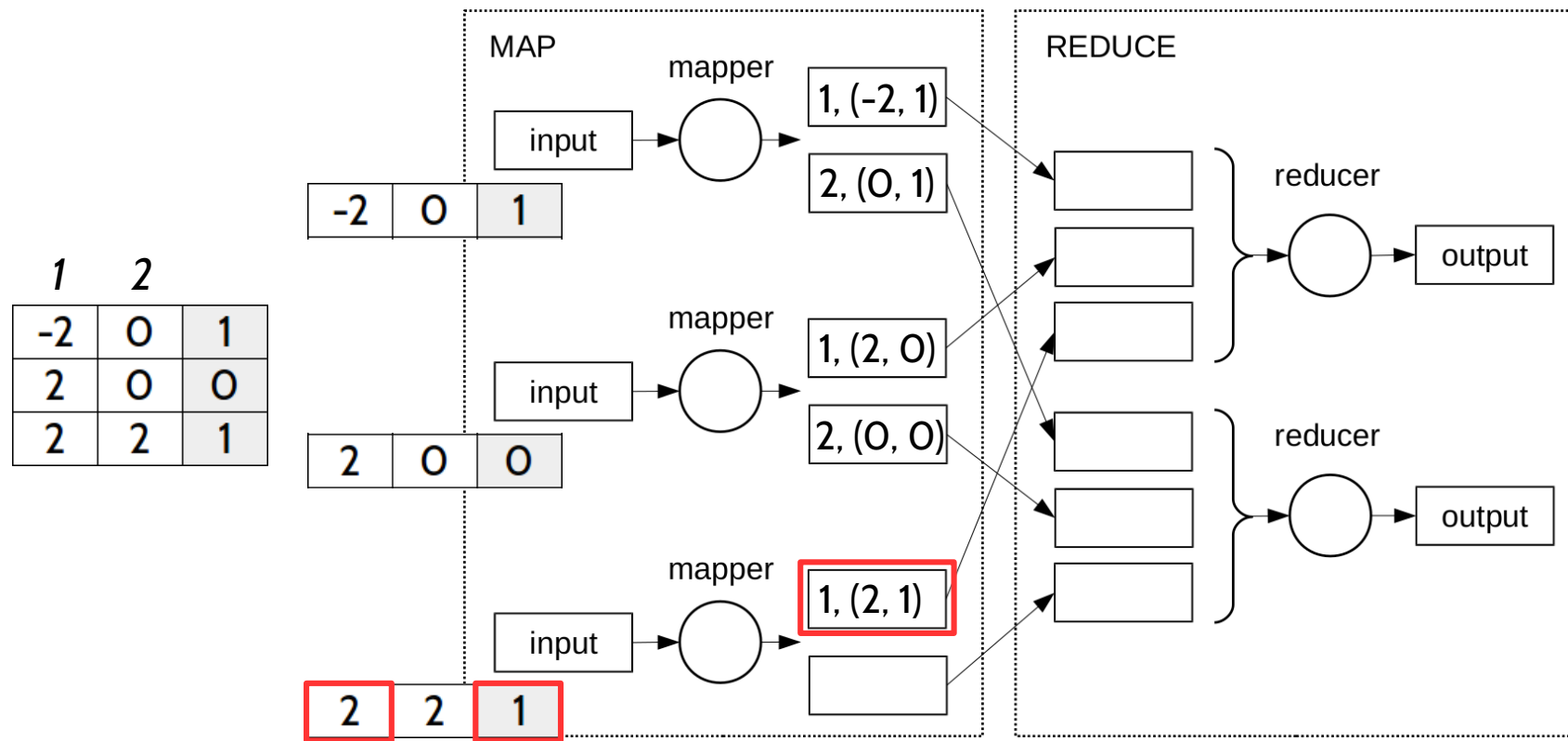# Brief introduction to MapReduce
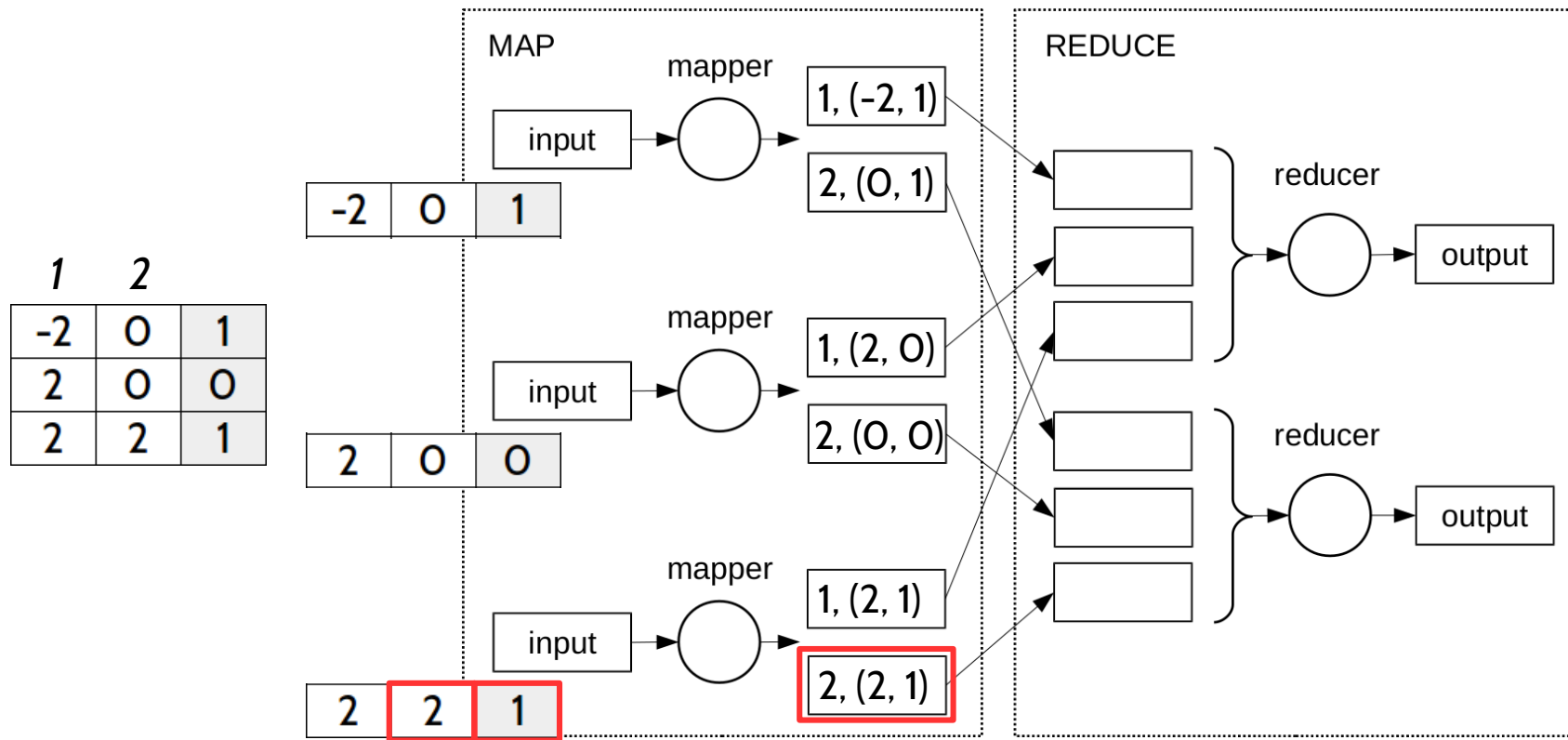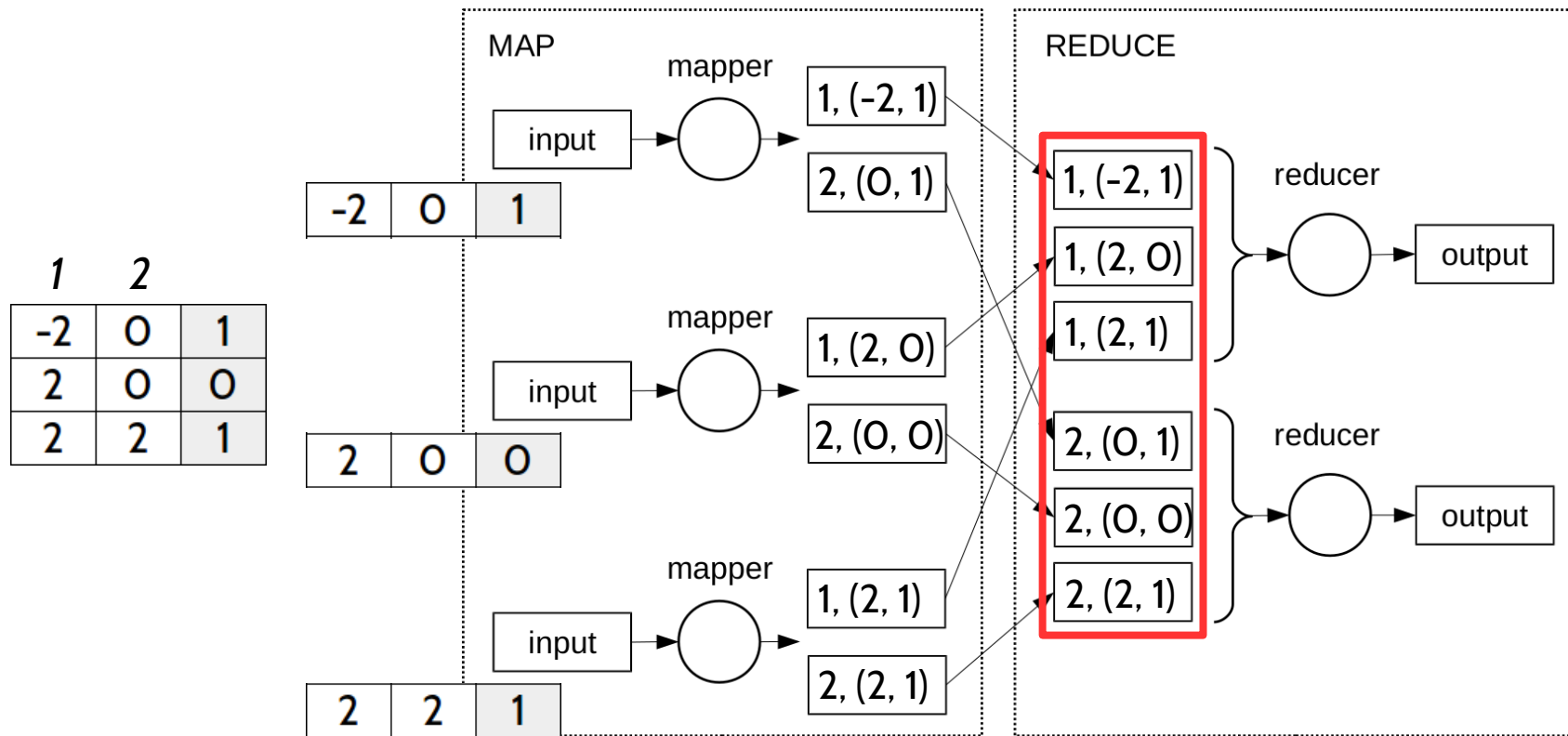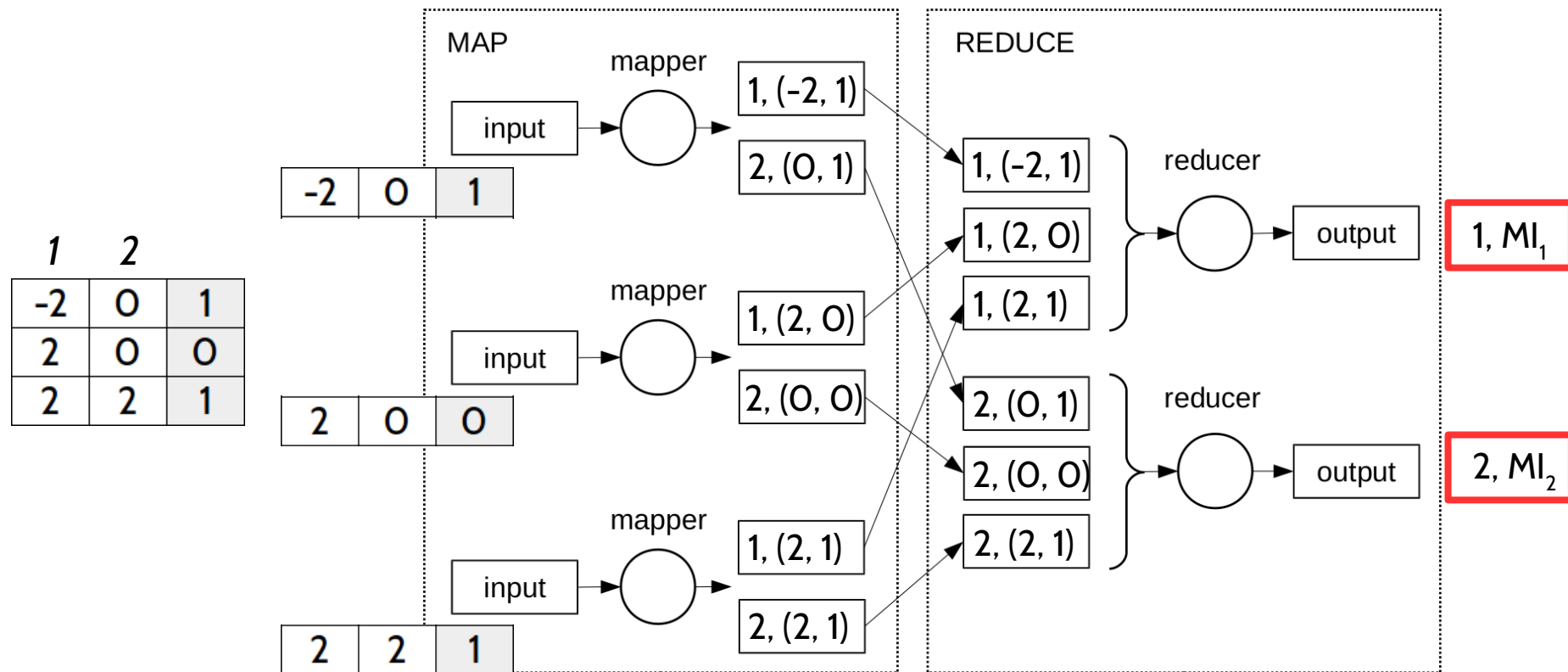
# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

| -2 | 0 | 1 |
|----|---|---|
| 2  | 0 | 0 |
| 2  | 2 | 1 |
| -2 | 0 | 1 |
| 2  | 0 | 0 |
| 2  | 2 | 1 |
| -2 | 0 | 1 |
| 2  | 0 | 0 |
| 2  | 2 | 1 |

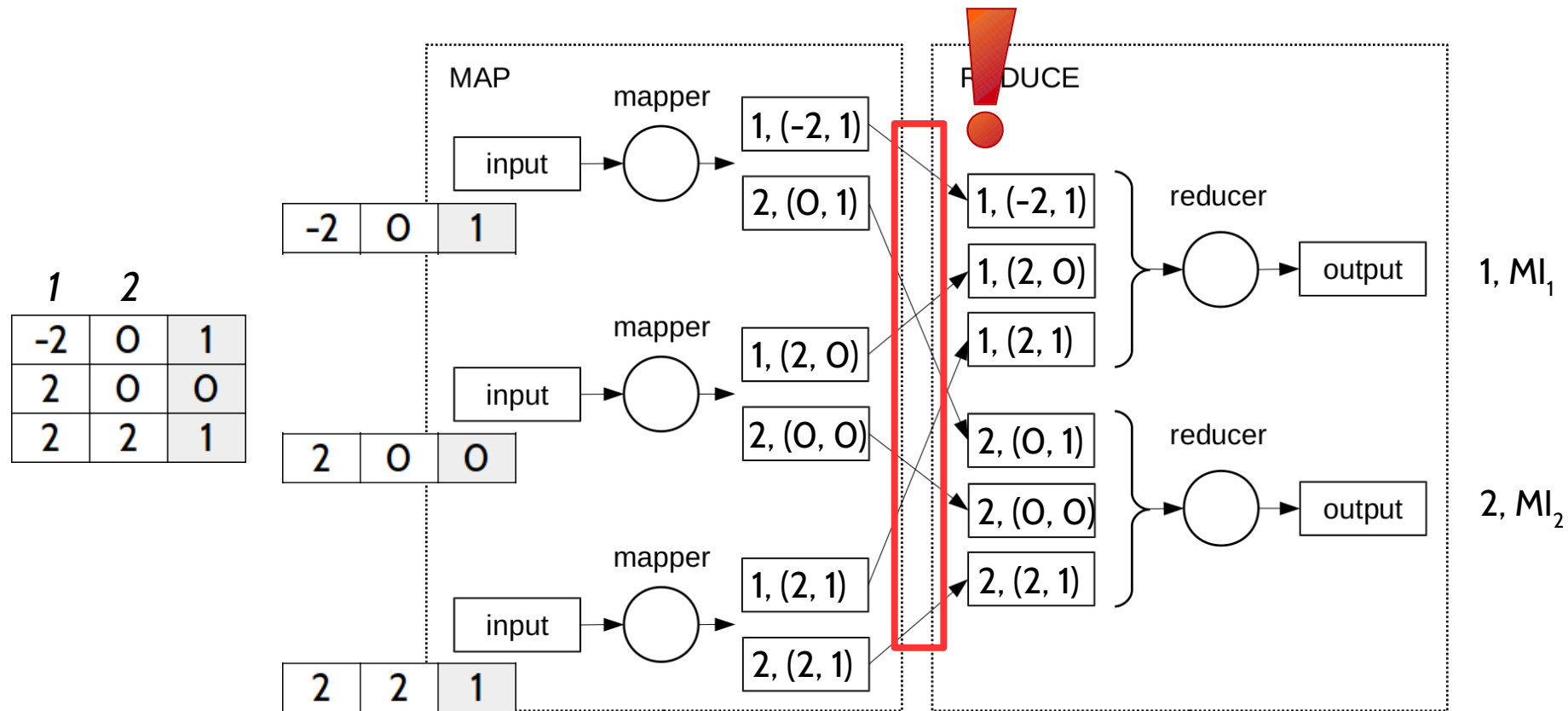...

| -2 | 0 | 1 |
|----|---|---|
| 2  | 0 | 0 |
| 2  | 2 | 1 |

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Brief introduction to MapReduce

# Data structure with cumulative property

# Data structure with cumulative property

MAP

| $x_1$ | c |
|-------|---|
| 2 | 0 |
| 0 | 0 |
| 0 | 0 |
| -2 | 1 |

# Data structure with cumulative property

MAP

|  | | $d_v$ | |
|---|---|---|---|
| | **-2** | **0** | **2** |

| $d_c$ | **0** |
|---|---|
| | **1** |

categories of the class

categories of the feature

| $\mathbf{x}_1$ | c |
|---|---|
| 2 | 0 |
| 0 | 0 |
| 0 | 0 |
| -2 | 1 |

$\longrightarrow$

# Data structure with cumulative property

MAP

|  |  | $d_v$ | | |
|---|---|---|---|---|
|  |  | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 0 | 1 |
|  | **1** | 0 | 0 | 0 |

| $\mathbf{x}_1$ | c |
|---|---|
| 2 | 0 |
| 0 | 0 |
| 0 | 0 |
| -2 | 1 |

$\longrightarrow$

# Data structure with cumulative property

MAP

|  | | $d_v$ | | |
|---|---|---|---|---|
|  | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 0 | 1 |
| | **1** | 0 | 0 | 0 |

|  | | $d_v$ | | |
|---|---|---|---|---|
|  | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 1 | 0 |
| | **1** | 0 | 0 | 0 |

| $\mathbf{x}_1$ | $\mathbf{c}$ |
|---|---|
| 2 | 0 |
| 0 | 0 |
| 0 | 0 |
| -2 | 1 |

# Data structure with cumulative property



MAP

| $\mathbf{x}_1$ | $\mathbf{c}$ |
|---|---|
| 2 | 0 |
| 0 | 0 |
| 0 | 0 |
| -2 | 1 |

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 0 | 1 |
| | **1** | 0 | 0 | 0 |

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 1 | 0 |
| | **1** | 0 | 0 | 0 |

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 1 | 0 |
| | **1** | 0 | 0 | 0 |

# Data structure with cumulative property

MAP

| $\mathbf{x}_1$ | $\mathbf{c}$ |
|---|---|
| 2 | 0 |
| 0 | 0 |
| 0 | 0 |
| -2 | 1 |

$\longrightarrow$

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 0 | 1 |
| | **1** | 0 | 0 | 0 |

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 1 | 0 |
| | **1** | 0 | 0 | 0 |

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 1 | 0 |
| | **1** | 0 | 0 | 0 |

| | | $d_v$ | | |
|---|---|---|---|---|
| | | **-2** | **0** | **2** |
| $d_c$ | **0** | 0 | 0 | 0 |
| | **1** | 1 | 0 | 0 |

# Data structure with cumulative property

# Data structure with cumulative property

# Data structure with cumulative property

# mRMR in MapReduce

# mRMR in MapReduce

1.    $i_c^1 = \{1, ..., N\}$

# mRMR in MapReduce

1   $i_c^1 = \{1, ..., N\}$
2   $i_s^1 = \varnothing$

# mRMR in MapReduce

$$1 \quad i_c^1 = \{1, ..., N\}$$
$$2 \quad i_s^1 = \varnothing$$
$$3 \quad \texttt{for} \quad l = 1 \rightarrow L$$

# mRMR in MapReduce

$$i_c^1 = \{1, ..., N\}$$

$$i_s^1 = \varnothing$$

```
for  l = 1 → L
    broadcast  i_class,  i_c^l,  i_s^l,  d_v,  d_c
    scores <- mapreduce(RDD, mapper, reducer)
```

# mRMR in MapReduce

$$i_c^1 = \{1, ..., N\}$$

$$i_s^1 = \varnothing$$

```
for l = 1 → L
    broadcast i_class, i_c^l, i_s^l, d_v, d_c
    scores <- mapreduce(RDD, mapper, reducer)
    k* ← collectArgmax(scores)
```

# mRMR in MapReduce



$$i_c^1 = \{1, ..., N\}$$
$$i_s^1 = \varnothing$$
```
for  l = 1 → L
   broadcast i_class,  i_c^l,  i_s^l,  d_v,  d_c
   scores <- mapreduce(RDD, mapper, reducer)
   k* ← collectArgmax(scores)
```
$$i_c^{l+1} \leftarrow i_c^l \setminus k^*$$
$$i_s^{l+1} \leftarrow i_s^l \cup k^*$$

# mRMR in MapReduce

$$i_c^1 = \{1, ..., N\}$$
$$i_s^1 = \varnothing$$
```
for  l = 1 → L
   broadcast i_class, i_c^l, i_s^l, d_v, d_c
   scores <- mapreduce(RDD, mapper, reducer)
   k* ← collectArgmax(scores)
```
$$i_c^{l+1} \leftarrow i_c^l \backslash k^*$$
$$i_s^{l+1} \leftarrow i_s^l \cup k^*$$
```
output  i_s^L
```

# mRMR in MapReduce, alternative layout

| $x_{1,1}$ | $x_{1,2}$ | $\ldots$ | $\ldots$ | $x_{1,N}$ |
|-----------|-----------|----------|----------|-----------|
| $x_{2,1}$ | $x_{2,2}$ | $\ldots$ | $\ldots$ | $x_{2,N}$ |
| $\ldots$  | $\ldots$  | $\ldots$ | $\ldots$ | $\ldots$  |
| $\ldots$  | $\ldots$  | $\ldots$ | $\ldots$ | $\ldots$  |
| $x_{M,1}$ | $x_{M,2}$ | $\ldots$ | $\ldots$ | $x_{M,N}$ |

# mRMR in MapReduce, alternative layout



| $x_{1,1}$ | $x_{1,2}$ | $\ldots$ | $\ldots$ | $x_{1,N}$ |
|---|---|---|---|---|
| $x_{2,1}$ | $x_{2,2}$ | $\ldots$ | $\ldots$ | $x_{2,N}$ |
| $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ |
| $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ |
| $x_{M,1}$ | $x_{M,2}$ | $\ldots$ | $\ldots$ | $x_{M,N}$ |

| $x_{1,1}$ | $x_{2,1}$ | $\ldots$ | $\ldots$ | $x_{M,1}$ |
|---|---|---|---|---|
| $x_{1,2}$ | $x_{2,2}$ | $\ldots$ | $\ldots$ | $x_{M,2}$ |
| $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ |
| $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ |
| $x_{1,N}$ | $x_{2,N}$ | $\ldots$ | $\ldots$ | $x_{M,N}$ |

# mRMR in MapReduce, alternative layout

# Scalability: #rows



mRMR with MI, ncol=1000, nfs=10, nexec=10

# Scalability: #columns



mRMR with MI, nrow=1M, nfs=10, nexec=10

# Scalability: #selected features



mRMR with MI, nrow=1M, ncol=50k, nexec=10

# Scalability: #nodes



mRMR with MI, nrow=1M, ncol=100, nfs=10

# Customization

$$\max_{x_j \in X - S_{m-1}} \left[ I(x_j; c) - \frac{1}{m-1} \sum_{x_i \in S_{m-1}} I(x_j; x_i) \right]$$

# Customization

$$\max_{x_j \in X - S_{m-1}} \left[ I(x_j; c) - \frac{1}{m-1} \sum_{x_i \in S_{m-1}} I(x_j; x_i) \right]$$

```
1  function getResult:
2    arguments:
3      variableArray: Array[Double]
4      classArray: Array[Double]
5      selectedVariablesArray:
            Array[Array[Double]]
6    return: Double
```

Yann-Aël
Le Borgne

Gianluca
Bontempi

github.com/creggian/spark-ifs

github.com/creggian/slides

claudioreggiani.com