# BigData: Apriori & SONs

A comparison among several implementations of Association-Rule based algorithms.

Cesare Montresor - VR481252

https://github.com/cesare-montresor/apriori-and-sons
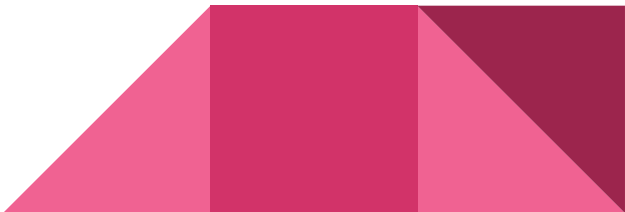
# Details

## Project

The project have been developed as framework for benchmarking a series of frequent itemset algorithms against various datasets and minimum support levels.
All the experiments and results have been obtained by trying as much as possible to replicate equal testing conditions. Loading of data, running of the experiments, profiling and result storage are unified and are not included in the computation of the timings, just the pure execution of the algorithm.

## Hardware

- Ubuntu 22.04
- 32GB RAM
- AMD Ryzen 9 5900X
- 12 Cores
- 24 Threads  *( 24 workers )*
- Spark 3.4.0 for Hadoop 3.3.4
- Python 3.10.9 ( conda )

## Algorithms

All implementations are based on PySpark and make use map, reduce, cartesian, filter and other typical paradigms from functional programming, replacing loops as much as possible.
Some implementations, like apriori_list are made in pure python as they are intended to be used on partition data, which by nature in Spark, is not an RDD, but simple python object. Also in this case the same guideline have been followed.
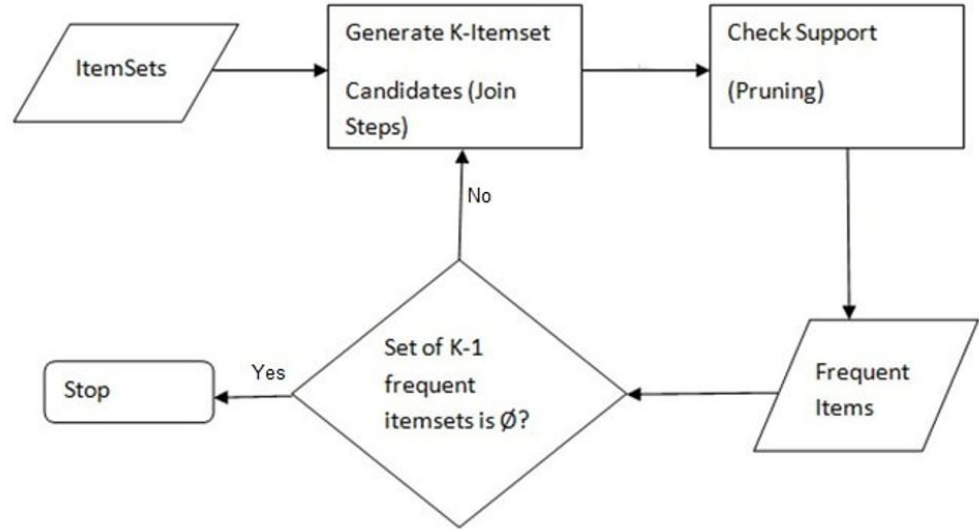
# Apriori

## Steps

1. Count single transaction items
2. Obtain FrequentSet of size 1
3. Permutate to obtain candidates size 2
4. Perform cartesian between candidates and transactions.
5. Filter only existing sets for existing transactions and count them.
6. Obtain FrequentSet on size 2
7. Repeat from step 3

# Apriori: Focus

## Extending the Set

```python
extendedSets = product(itemSupport, itemUnique) ⭐
extendedSets = map(flatMergeCartesian, extendedSets)  # ( ('3','1'),'4') -> ('1','3','4')
extendedSets = tuple(filter(lambda item: len(item) == set_length, extendedSets))
# from flatMergeCartesian: ( ('4','3'), '3') -> ('3','4')
extendedSets = tuple(set(extendedSets))
```
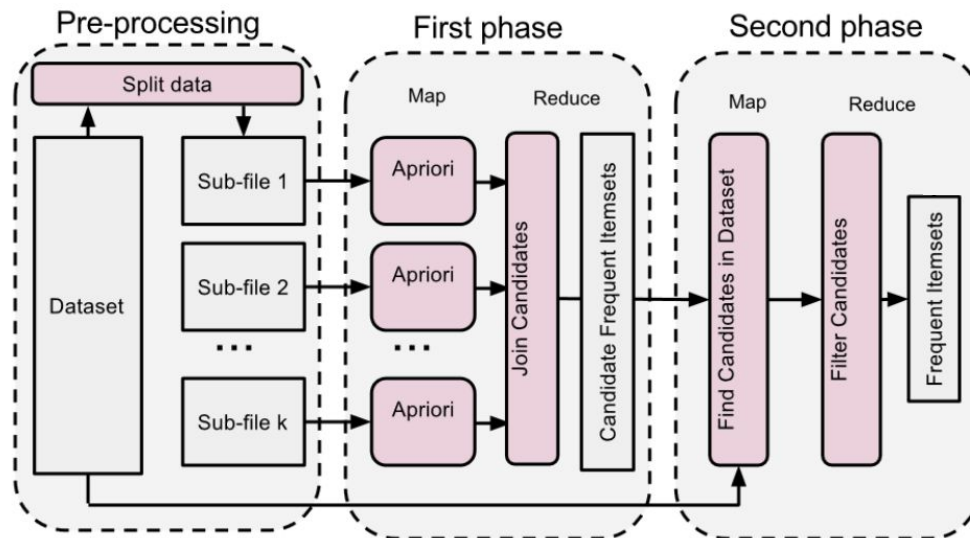
## Filtering the Set

```python
existingSets = product(extendedSets, transactions) ⭐
existingSets = tuple(filter(lambda item: set(item[0]).issubset(item[1]), existingSets))
existingSets = tuple(map(lambda item: item[0], existingSets))
```

# SON

## Steps

1. Partition the dataset
2. Perform "apriori" for each partition, discard the counting.
3. Merge the results from all partitions and use those as candidates.
4. Perform cartesian between candidates and transactions.
5. Find candidate sets.
6. Filter candidates and obtain Frequent Itemsets.

1

2

# SON: Focus

Find candidates from partitions

```
candidateSets = partitionSets.mapPartitions(lambda part: ⭐apriori_list(list(part), min_support))
candidateSets = candidateSets.map(lambda item: item[0])
candidateSets = candidateSets.coalesce(1)
candidateSets = candidateSets.distinct()
```
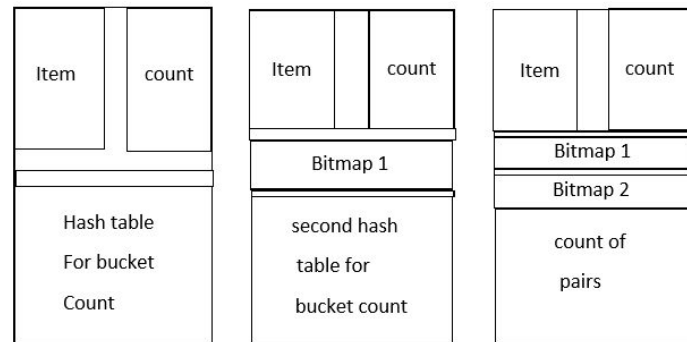
Test candidates

```
# test candidates
testingSets = transactions.cartesian(candidateSets) ⭐
testingSets = testingSets.filter(lambda item: set(item[0]).issubset(item[1]))
testingSets = testingSets.reduceByKey(lambda x, y: x + y)
testingSets = testingSets.filter(lambda item: item[1] > candidate_frequency)
```

# PCY

## Steps

1. Count single transaction items
2. Obtain FrequentSet of size 1
3. Partition the dataset into buckets
4. Generate Candidate set from row combination.
5. Using a hash function count items per bucket and extract candidates and use it for filtering.
6. Perform cartesian between candidates and transactions.
7. Find candidate sets.
8. Filter candidates and obtain Frequent Itemsets.



Multistage Memory Mapping.

# PCY: Key ideas

Each row contains, in potential, all the frequent/candidate itemsets.

Make efficient use of memory it's important, use hash tables to select candidates.

Frequent Itemsets of size N, are a combination, without repetition of the individual items of N-1.

# Why PCY?

**?**

Each row can be exploded in it subsets: processing can be parallelized at row level.

Why buckets ?

**?**

You can sum and obtain final frequency.

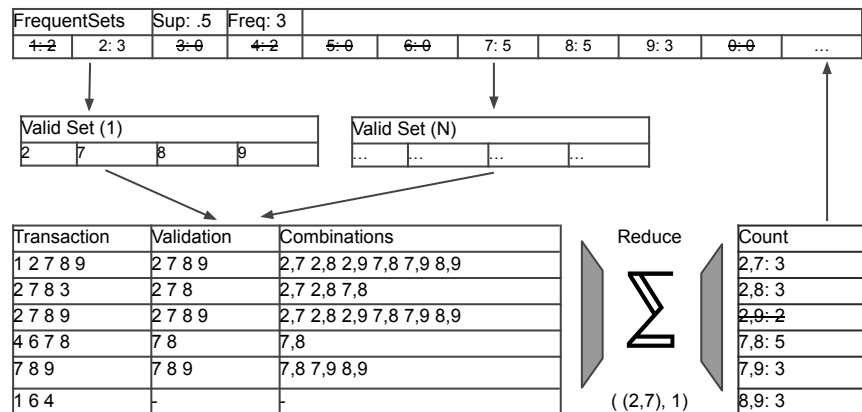Why multiply for hashing, to generate just candidates?

**?**

Remove every item you can before join and combinations.

Why joining huge tables?

# PCY Bucketless

## Steps

1. Count single transaction items
2. Obtain FrequentSets of size N=1.
3. Make the valid set, the FrequentSets on size N.
4. For each valid transaction line generate a list of actual sets of size N from the combination items of that row.
5. Each transaction line generates multiple sets, make them flat.
6. Reduce by key and filter.
7. Obtain Frequent Itemsets.

| FrequentSets | Sup: .5 | Freq: 3 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1: 2 | 2: 3 | 3: 0 | 4: 2 | 5: 0 | 6: 0 | 7: 5 | 8: 5 | 9: 3 | 0: 0 | ... |

| Valid Set (1) | | | |
|---|---|---|---|
| 2 | 7 | 8 | 9 |

| Valid Set (N) | | | |
|---|---|---|---|
| ... | ... | ... | ... |

| Transaction | Validation | Combinations |
|---|---|---|
| 1 2 7 8 9 | 2 7 8 9 | 2,7 2,8 2,9 7,8 7,9 8,9 |
| 2 7 8 3 | 2 7 8 | 2,7 2,8 7,8 |
| 2 7 8 9 | 2 7 8 9 | 2,7 2,8 2,9 7,8 7,9 8,9 |
| 4 6 7 8 | 7 8 | 7,8 |
| 7 8 9 | 7 8 9 | 7,8 7,9 8,9 |
| 1 6 4 | - | - |

Reduce

$$\sum$$

( (2,7), 1)

| Count |
|---|
| 2,7: 3 |
| 2,8: 3 |
| 2,9: 2 |
| 7,8: 5 |
| 7,9: 3 |
| 8,9: 3 |

# PCY Bucketless: Focus

## Reduce Transactions

```
pairsCount = transactions.flatMap(lambda transaction: findPairs(transaction, validationSet, candidateSize))
pairsCount = pairsCount.reduceByKey(lambda x, y: x + y)
pairsCount = pairsCount.filter(lambda item: item[1] > min_frequency)
```

## Find Pairs

```
def findPairs(transaction:list, validationSet:tuple, candidateSize:int):
    if len(transaction) < candidateSize: return None
    validItems = set(filter(lambda item: item in validationSet, transaction))
    if len(validItems) < candidateSize: return tuple()

    pairs = combinations(transaction, candidateSize)
    return tuple((pair, 1) for pair in pairs)
```

# How far can I push it?

Why not consolidating the item removal? It will avoid redundant future checks.
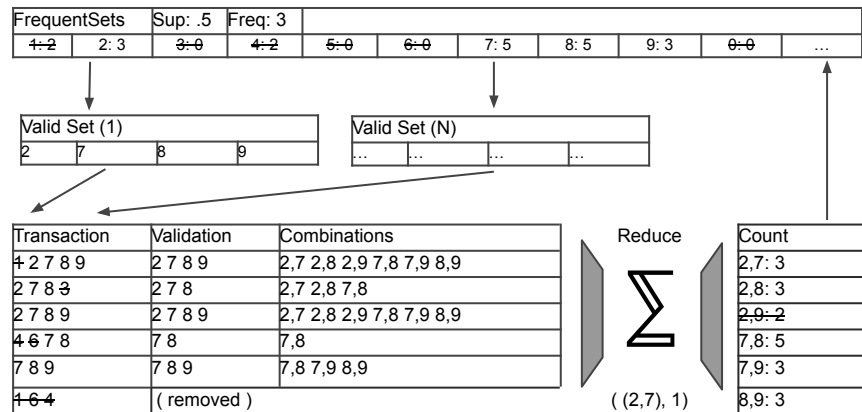
Keep calm and just count frequencies.

Progressive expansion of N allow to cheaply apply Erosion keeping the future expansions limited.

# PCY Erode

## Steps

1. Count single transaction items
2. Obtain FrequentSets of size N=1.
3. Make the ValidationSet, the FrequentSets on size N.
4. Remove each line with less then N items.
5. Remove from each transaction line the items not present in the itemsets.
6. Remove each line with less then N item.
7. For each transaction line generate a list of actual sets of size N from the combination items of that row.
8. Each transaction line generates multiple sets, make them flat.
9. Reduce by key and filter.
10. Filter candidates and obtain Frequent Itemsets.

| FrequentSets | Sup: .5 | Freq: 3 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| ~~1: 2~~ | 2: 3 | ~~3: 0~~ | ~~4: 2~~ | ~~5: 0~~ | ~~6: 0~~ | 7: 5 | 8: 5 | 9: 3 | ~~0: 0~~ | ... |

| Valid Set (1) | | | |
|---|---|---|---|
| 2 | 7 | 8 | 9 |

| Valid Set (N) | | | |
|---|---|---|---|
| ... | ... | ... | ... |

| Transaction | Validation | Combinations |
|---|---|---|
| ~~1~~ 2 7 8 9 | 2 7 8 9 | 2,7 2,8 2,9 7,8 7,9 8,9 |
| 2 7 8 ~~3~~ | 2 7 8 | 2,7 2,8 7,8 |
| 2 7 8 9 | 2 7 8 9 | 2,7 2,8 2,9 7,8 7,9 8,9 |
| ~~4 6~~ 7 8 | 7 8 | 7,8 |
| 7 8 9 | 7 8 9 | 7,8 7,9 8,9 |
| ~~1 6 4~~ | ( removed ) | |

Reduce

$$\sum$$

( (2,7), 1)

| Count |
|---|
| 2,7: 3 |
| 2,8: 3 |
| ~~2,9: 2~~ |
| 7,8: 5 |
| 7,9: 3 |
| 8,9: 3 |

# PCY Erode: Focus

## Reduce Transactions

```python
transactions = reduceTransactions(transactions, validationSet, candidateSize)

def reduceTransactions(transactions: pyspark.RDD, validationSet, candidateSize:int)->pyspark.RDD:
    reducedSet = transactions.map(lambda transaction: filterTransactions(transaction, validationSet, candidateSize))
    reducedSet = reducedSet.filter(lambda transaction: transaction is not None)
    return reducedSet

def filterTransactions(transaction, validationSet, candidateSize:int):
    if len(transaction) < candidateSize: return None
    validItems = tuple(filter(lambda item: item in validationSet, transaction))
    if len(validItems) < candidateSize: return None
    return validItems
```
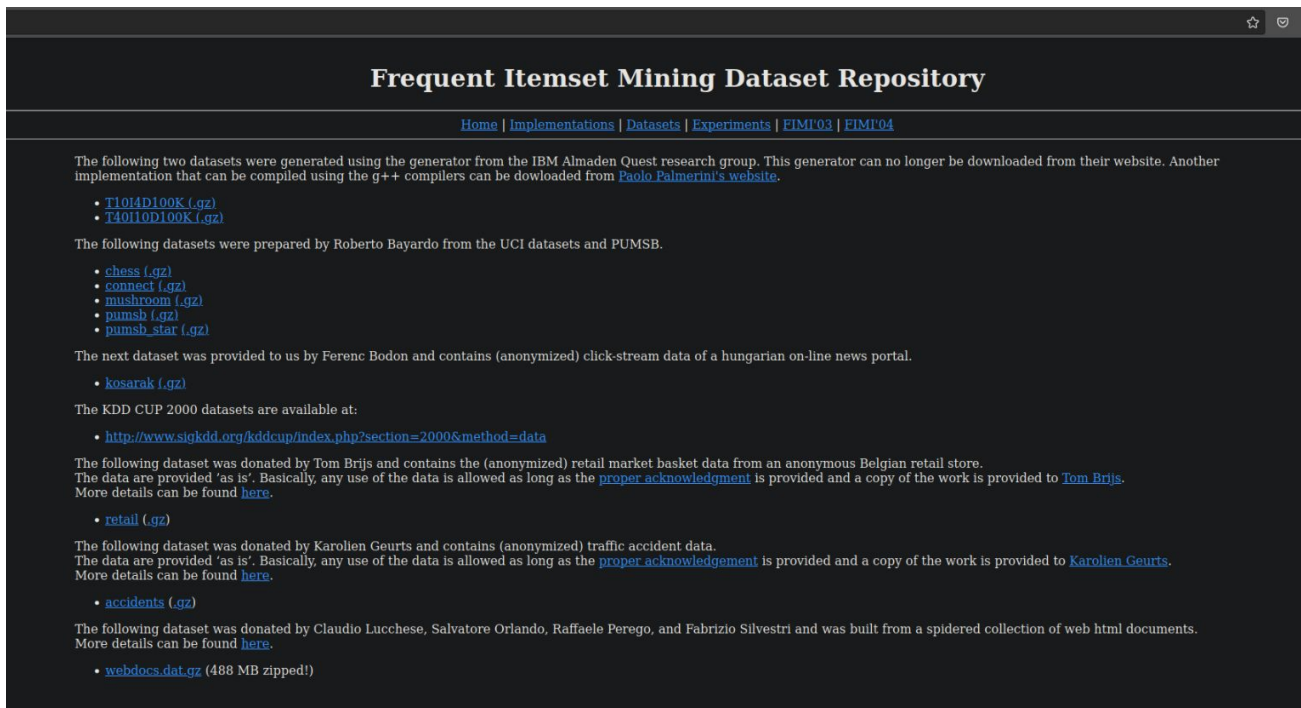
## Find Pairs

```python
def findPairs(transaction:list|tuple, candidateSize:int) -> tuple:
    pairs = combinations(transaction, candidateSize)
    return tuple((pair, 1) for pair in pairs)
```

# DATASETS

# Dataset: Overview

# Dataset: Overview

| Name | Size | Lines | Num Items | Support | | AVG Items |
|------|------|-------|-----------|---------|-----|-----------|
| | | | | AVG | MAX | |
| Mushroom | 570.408 | 8.124 | 119 | 0.193 | 1.000 | 23 |
| Retail | 4.167.490 | 88.162 | 16.470 | 0.001 | 0.575 | 10 |
| 50K10 | 2.008.453 | 50.000 | 869 | 0.012 | 0.078 | 10 |
| 25K40 | 3.869.620 | 25.000 | 942 | 0.042 | 0.288 | 40 |
| 100K10 | 4.022.055 | 100.000 | 870 | 0.0116 | 0.078 | 10 |
| 100K40 | 15.478.113 | 100.000 | 942 | 0.042 | 0.287 | 40 |

# Dataset: Mushroom

## General

| Size | 570.408 |
|---|---|
| Num lines | 8.124 |
| Num Items | 119 |
| Support avg | 0.193 |
| Support max | 1 |

## Frequency

| Freq min | 4 |
|---|---|
| Freq max | 8.124 |
| Freq avg | 1.570 |

## Transactions

| Items min | 23 |
|---|---|
| items max | 23 |
| items avg | 23 |

### Items frequencies



### Items per transaction

# Dataset: Retail

## General

| Size | 4.167.490 |
|---|---|
| Num lines | 88.162 |
| Num Items | 16.470 |
| Support avg | 0.001 |
| Support max | 0.575 |

## Frequency

| Freq min | 1 |
|---|---|
| Freq max | 50.675 |
| Freq avg | 55 |

## Transactions

| Items min | 1 |
|---|---|
| items max | 76 |
| items avg | 10 |

### Items frequencies



### Items per transaction

# Dataset: 50K10

### General

| Size | 2.008.453 |
|---|---|
| Num lines | 50.000 |
| Num Items | 869 |
| Support avg | 0.012 |
| Support max | 0.078 |

### Frequency

| Freq min | 2 |
|---|---|
| Freq max | 3.905 |
| Freq avg | 581 |

### Transactions

| Items min | 1 |
|---|---|
| items max | 29 |
| items avg | 10 |

**Items frequencies**



**Items per transaction**

# Dataset: 25K40

### General

| Size | 3.869.620 |
|---|---|
| Num lines | 25.000 |
| Num Items | 942 |
| Support avg | 0.042 |
| Support max | 0.288 |

### Frequency

| Freq min | 1 |
|---|---|
| Freq max | 7188 |
| Freq avg | 1051 |

### Transactions

| Items min | 9 |
|---|---|
| items max | 72 |
| items avg | 40 |

## Items frequencies



## Items per transaction

# Dataset: 100K10

### General

| Size | 4.022.055 |
|---|---|
| Num lines | 100.000 |
| Num Items | 870 |
| Support avg | 0.012 |
| Support max | 0.078 |

### Frequency

| Freq min | 1 |
|---|---|
| Freq max | 7828 |
| Freq avg | 1161 |

### Transactions

| Items min | 1 |
|---|---|
| items max | 29 |
| items avg | 10 |

### Items frequencies



### Items per transaction

# Dataset: 100K40

## General

| Size | 15.478.113 |
|---|---|
| Num lines | 100.000 |
| Num Items | 942 |
| Support avg | 0.042 |
| Support max | 0.287 |

## Frequency

| Freq min | 5 |
|---|---|
| Freq max | 28738 |
| Freq avg | 4204 |

## Transactions

| Items min | 4 |
|---|---|
| items max | 77 |
| items avg | 40 |

### Items frequencies



### Items per transaction

# RESULTS

# Results:

| | Support | Itemset | apriori | | pcy | pcy erode | son |
|---|---|---|---|---|---|---|---|
| | | | py | spark | | | |
| mushroom | 0.9 | **9** | **6** | **13** | **6** | **3** | **8** |
| | 0.7 | **31** | **18** | **111** | **164** | **4** | **9** |
| | 0.5 | **153** | **81** | **136** | **1.417** | **5** | **16** |
| | 0.3 | **2735** | **1.423** | OOM | **24.399** | **93** | **140** |

| | Support | Itemset | py | spark | pcy | pcy erode | son |
|---|---|---|---|---|---|---|---|
| retail | 0.5 | **1** | **531** | **96** | **3** | **3** | **60** |
| | 0.3 | **3** | **1.577** | **253** | **7** | **3** | **98** |
| | 0.1 | **9** | **4.803** | **723** | **140** | **3** | **221** |

| | Support | Itemset | py | spark | pcy | pcy erode | son |
|---|---|---|---|---|---|---|---|
| 25K40 | 0.1 | **81** | **1.740** | **254** | **9** | **3** | **137** |
| 100K40 | 0.1 | **82** | **7.115** | **1.040** | **26** | **6** | **940** |

| | Support | Itemset | py | spark | pcy | pcy erode | son |
|---|---|---|---|---|---|---|---|
| 50K10 | 0.01 | **384** | - - - | - - - | - - - | **4** | **426** |
| | 0.03 | **58** | - - - | - - - | - - - | **2** | **94** |
| | 0.06 | **4** | - - - | - - - | - - - | **2** | **31** |
| retail | 0.005 | **580** | - - - | - - - | - - - | **7** | **12.298** |
| | 0.05 | **16** | - - - | - - - | - - - | **3** | **339** |

# Results: Cost of a zero

| | | apriori | | pcy | pcy erode | son |
|---|---|---|---|---|---|---|
| | | py | spark | | | |
| retail | 0.7 | **1** | **2** | **3** | **2** | **41** |
| | 0.9 | **1** | **2** | **2** | **2** | **41** |
| 100K10 | 0.1 | **1** | **2** | **2** | **2** | **46** |
| | 0.3 | **1** | **2** | **2** | **2** | **45** |
| | 0.5 | **1** | **2** | **3** | **2** | **44** |
| | 0.7 | **1** | **2** | **2** | **2** | **45** |
| | 0.9 | **1** | **2** | **2** | **2** | **44** |
| 100K40 | 0.3 | **2** | **3** | **3** | **3** | **45** |
| | 0.5 | **2** | **3** | **3** | **3** | **45** |
| | 0.7 | **2** | **3** | **3** | **3** | **45** |
| | 0.9 | **2** | **3** | **3** | **3** | **45** |

| | | apriori | | pcy | pcy erode | son |
|---|---|---|---|---|---|---|
| | | py | spark | | | |
| 50K10 | 0.1 | **1** | **2** | **2** | **2** | **24** |
| | 0.3 | **1** | **2** | **2** | **2** | **24** |
| | 0.3 | **1** | **2** | **2** | **2** | **24** |
| | 0.7 | **1** | **2** | **2** | **2** | **25** |
| | 0.9 | **1** | **2** | **2** | **2** | **25** |
| 100K40 | 0.3 | **2** | **3** | **3** | **3** | **45** |
| | 0.5 | **2** | **3** | **3** | **3** | **45** |
| | 0.7 | **2** | **3** | **3** | **3** | **45** |
| | 0.9 | **2** | **3** | **3** | **3** | **45** |

# Erosion at work: 50K10 @ 0.01 (500)

```
========================================================================
                      pcy_erode  -50K10-0010
========================================================================
Min Freq:  500 = 50000 * 0.01
374
candidateSize:  2
supportSets: ['240', '274', '538', '630', '825', '834', '581', '814', '674', '733', '854', '950', '422', '449', '857', '229', '283', '738', '853', '883', '966', '978', '143', '185', '214', '658', '682',
'782', '947', '970', '227', '390', '192', '208', '279', '280', '496', '530', '675', '720', '914', '932', '217', '161', '175', '490', '571', '623', '960', '130', '461', '862', '900', '147', '411', '572',
'579', '778', '290', '458', '70', '204', '334', '513', '504', '73', '419', '469', '722', '846', '326', '526', '975', '116', '198', '541', '805', '631', '780', '935', '17', '763', '956', '145', '385', '676',
'790', '792', '885', '522', '12', '296', '354', '548', '346', '477', '829', '234', '649', '600', '157', '115', '517', '736', '744', '641', '417', '628', '111', '154', '580', '10', '132', '21', '54', '348',
'100', '48', '319', '112', '140', '285', '387', '93', '583', '122', '718', '1', '69', '797', '110', '509', '611', '995', '33', '336', '598', '470', '992', '897', '259', '45', '162', '378', '716', '8', '413',
'823', '982', '515', '694', '57', '812', '414', '752', '361', '108', '486', '440', '265', '540', '468', '819', '886', '429', '68', '4', '887', '707', '815', '948', '634', '351', '949', '163', '335', '922',
'173', '258', '608', '820', '207', '25', '52', '368', '448', '561', '687', '775', '39', '120', '205', '401', '704', '35', '895', '937', '964', '294', '381', '708', '766', '104', '569', '620', '798', '350',
'529', '809', '71', '597', '618', '183', '276', '653', '706', '878', '177', '424', '795', '910', '125', '392', '27', '78', '921', '803', '266', '523', '614', '888', '944', '43', '480', '874', '151', '890',
'310', '810', '844', '918', '967', '403', '774', '788', '789', '201', '171', '701', '946', '471', '487', '638', '678', '735', '242', '758', '617', '859', '684', '740', '841', '210', '605', '884', '460',
'746', '28', '5', '919', '196', '489', '494', '673', '362', '591', '31', '58', '181', '472', '573', '651', '168', '632', '832', '871', '988', '72', '981', '32', '239', '500', '126', '639', '765', '521',
'594', '606', '236', '952', '90', '593', '941', '423', '516', '6', '913', '577', '343', '527', '989', '97', '574', '793', '427', '37', '55', '275', '51', '534', '906', '576', '373', '665', '963', '349',
'197', '749', '94', '984', '692', '567', '800', '41', '923', '377', '991', '998', '899', '710', '867', '170', '438', '563', '357', '332', '322', '928', '75', '38', '784', '686', '663', '843', '129', '578',
'510', '860', '309', '804', '826', '394', '105', '308', '661', '405', '688', '893', '85', '450', '550', '769', '554', '366']
candidateSize:  3
supportSets: ['227', '390', '217', '346', '722', '825', '829', '682', '39', '704', '789', '368']
candidateSize:  4
supportSets: ['825', '39', '704']
------------- RESULTS -------------
pcy_erode -50K10-0010
success
384
3.919753313064575
------------------------------------
TASK: success  384  3.919753313064575
```

# Erosion at work: retail @ 0.005 (440)

```
================================================================
                    pcy_erode  -retail-0005
================================================================
Min Freq: 440 = 88162 * 0.005
221
candidateSize: 2
supportSets: ['9', '10', '19', '45', '48', '53', '56', '60', '107', '110', '147', '150', '161', '175', '179', '185', '208', '209', '229', '237', '249', '255', '258', '259', '264', '270', '272', '279', '286',
'301', '334', '338', '345', '365', '389', '408', '413', '426', '441', '449', '464', '488', '490', '522', '548', '570', '571', '581', '623', '649', '664', '675', '677', '703', '772', '783', '790', '806',
'812', '824', '831', '846', '885', '947', '956', '1020', '1043', '1144', '1146', '1198', '1239', '1344', '1404', '1600', '1659', '1693', '1704', '1783', '1814', '1872', '2046', '2080', '2118', '2168', '2383',
'2879', '2958', '4883', '10515', '12929', '15618', '16011', '2', '11', '18', '23', '30', '31', '32', '36', '37', '38', '39', '41', '47', '49', '52', '55', '62', '65', '76', '78', '79', '80', '89', '94',
'101', '103', '105', '117', '123', '136', '155', '156', '170', '178', '186', '201', '225', '242', '251', '260', '261', '267', '269', '271', '281', '297', '310', '340', '371', '374', '379', '396', '398',
'405', '407', '418', '420', '425', '432', '438', '475', '479', '533', '535', '544', '549', '589', '592', '604', '681', '704', '740', '766', '789', '793', '798', '808', '809', '855', '856', '865', '878',
'910', '913', '916', '961', '976', '1004', '1067', '1113', '1135', '1327', '1355', '1393', '1417', '1435', '1479', '1578', '1585', '1677', '1714', '1715', '1867', '1986', '1987', '2051', '2135', '2184',
'2199', '2238', '2284', '2329', '2399', '2437', '3270', '3616', '4336', '8978', '10444', '10446', '12925', '12946', '12959', '13041', '14098', '14933', '15832', '16010', '16217']
candidateSize: 3
supportSets: ['147', '48', '175', '179', '229', '249', '255', '258', '237', '286', '161', '334', '338', '60', '413', '45', '270', '522', '107', '264', '548', '570', '110', '19', '649', '783', '812', '824',
'9', '956', '1146', '677', '301', '1600', '703', '1814', '185', '772', '2958', '10515', '56', '272', '664', '790', '806', '10', '389', '1344', '12929', '16011', '36', '38', '39', '41', '37', '55', '49', '32',
'76', '78', '79', '101', '89', '103', '105', '117', '170', '201', '65', '242', '225', '269', '271', '310', '371', '11', '405', '438', '475', '549', '592', '604', '23', '533', '740', '789', '479', '31', '589',
'766', '123', '1135', '1327', '1393', '1435', '1004', '1578', '544', '18', '976', '2238', '156', '3270', '13041', '12925', '14098', '15832', '16010', '16217']
candidateSize: 4
supportSets: ['147', '48', '179', '249', '255', '258', '237', '286', '60', '270', '548', '110', '824', '9', '1146', '677', '301', '338', '2958', '185', '413', '36', '38', '39', '41', '37', '89', '170', '32',
'225', '105', '371', '65', '475', '271', '310', '2238', '49', '79', '101', '78', '438', '592', '201', '589', '533', '123', '604', '1327', '1393', '13041', '12925', '14098', '16010', '16217']
candidateSize: 5
supportSets: ['286', '48', '110', '36', '38', '39', '41', '170', '32', '89', '65', '225', '310']
candidateSize: 6
supportSets: ['48', '32', '38', '39', '41']
------------- RESULTS -------------
pcy_erode-retail-0005
success
580
6.335402727127075
-----------------------------------
TASK: success  580  6.335402727127075
```

# Erosion at work: 25K40 @ 0.01 (250)

========================================================
                    pcy_erode-25K40-0010
========================================================
Min Freq: 250 = 25000 * 0.01
750
candidateSize: 2
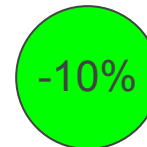
# Other findings

## Hadoop native libraries

```
23/07/02 14:44:47 WARN NativeCodeLoader:
Unable to load native-hadoop library for your platform... using builtin-java classes where applicable.
```

NO GAIN

## Using Integers

```python
def load(sc, data_path, sep=' ', useInt=True) -> pyspark.RDD:
    itemSets = sc.textFile(data_path)
    itemSets = itemSets.map(lambda line: line.strip().split(sep))
    if useInt:
        itemSets = itemSets.map(lambda line: tuple(map(lambda item: int(item), line)))
    return itemSets
```

-10%

| | | Sup | Strings | Integers | Time saved | |
|---|---|---|---|---|---|---|
| mushroom | apriori spark | 0.7 | 113 | 115 | -0% | ← *(pure pyspark)* |
| | pcy erode | 0.3 | 93 | 77 | -15% | |
| | apriori python | 0.5 | 81 | 69 | -10% | |

# CONCLUSIONS

# Conclusions

- <u>Always plot stats of your dataset as first step, will save tons of time down the line.</u>
- Curse of dimensionality:
  - Join and Catensian products should be avoid or limited.
  - Alternatives: in-transaction line combinations
  - *Better a Join then a loop: combinations can be an option but can become triky.*
- Filtering can help reducing the overall computation and if done cheaply doesn't decrease the performance in the worst case scenario.
- Performing an initial filtering even against the Frequent Itemset of size 1, which is anyway computed, can significantly reduce the number items per line, as well as overall lines, greatly decreasing.
- Finding and Verifying Candidate Itemsets can be both wasteful and very expensive.
- In-line combinations when used in combination with counting and filtering can be a scalable alternative to often wasteful cartesian products.
- PCY Erode shows that progressive erosion of the dataset and controlled expansions of the rows, an effective ways to keep under control faster-then-linear growths.
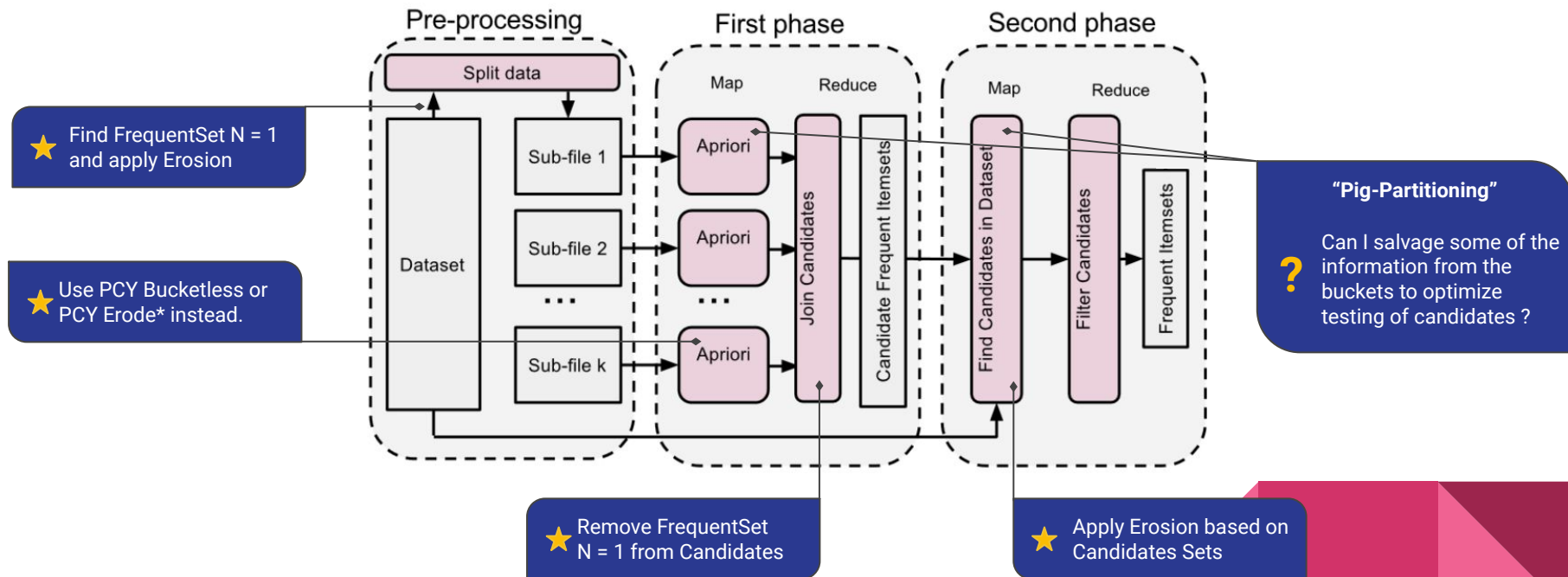
FUTURE

# Future

- By the current results I foresee a potential for proving "PCY Erode" with even more precise line filtering and while removing the need of resorting to combinations, potentially reaching an almost guaranteed minimum expansions of each row. However the computational cost of implementing such an aggressive filtering has yet to be accessed and trade-offs has yet to be evaluated.

- Considering that "PCY Erode" doesn't use the concept of buckets should be a good algorithm to replace apriori to discover Candidate Itemsets.

- By the current results I believe it should be possible to apply several optimizations also to the second phase of the SON algorithm by reusing the partitioning generated at the first step, combined with filtering and in-line expansion, to sensibly reduce the size of the cartesian product between the transaction and the candidate itemset.

# Future: SON Erode



Pre-processing

Split data

Find FrequentSet N = 1 and apply Erosion

Use PCY Bucketless or PCY Erode* instead.

Dataset

Sub-file 1

Sub-file 2

...

Sub-file k

First phase

Map    Reduce

Apriori

Apriori

...

Apriori

Join Candidates

Candidate Frequent Itemsets

Second phase

Map    Reduce

Find Candidates in Dataset

Filter Candidates

Frequent Itemsets

**"Pig-Partitioning"**

? Can I salvage some of the information from the buckets to optimize testing of candidates ?

Remove FrequentSet N = 1 from Candidates

Apply Erosion based on Candidates Sets

* need to make a copy of the partition