

Mapreduce – Duplikate finden und auflisten

```
from mrsim import mr simulator
def map(key, val):
    # res = List
    res = []
    # (key, val) = Tupel, der der Liste hinzugefügt wird
    res.append((key, val))
    return res
def reduce(key, val):
    res = []
    if len(val) > 1:
        # removeduplicates (original order is not preserved)
        val = list(set(val))
        # res.append((key, val))
        res.append((str(key) + ':' + str(val)))
    return res
inputs = [(123, 'Name1'), (123, 'Name2'), (456, 'Name1'),
          (345, 'Name3'), (456, 'Name2'), (123, 'Name1')]
res = mr simulator(inputs, map, reduce)
print(res)
```

Python Suchen

a) Merge Sort:

```
#merge function
def merge(arr1, arr2):
    result = []
    while len(arr1) != 0 and len(arr2) != 0:
        if arr1[0] < arr2[0]:
            result.append(arr1[0])
            arr1.remove(arr1[0])
        else:
            result.append(arr2[0])
            arr2.remove(arr2[0])
    if len(arr1) == 0:
        result += arr2
    else:
        result += arr1
    return result
#Code for mergesort
def mergesort(arr):
    if len(arr) <= 1:
        return arr
    middle = len(arr) // 2
    left = mergesort(arr[:middle])
    right = mergesort(arr[middle:])
    return merge(left, right)
#test
print(mergesort([3, 6, 8, 10, 1, 2, 1]))
1
```

b) Depth-first Search:

```
def deepSearch(graph, start, destination, visitedNodes=[]):
    if start == destination:
        return True
    if graph[start]:
        for node in filter(lambda x: x not in visitedNodes,
                           graph[start]):
            visitedNodes.append(node)
    if deepSearch(graph, node, destination, visitedNodes):
        return True
```

```

else:
    visitedNodes.append('Node'+str(node)+
        'is finished, no path from'+str(node)+
        'to'+str(destination))
    return False
adjList=f0:[1,2],1:[2,3],2:[4],3:[4,5],4:[],5:[] g
visitedNodes=[]
print(deepSearch(adjList,0,5,visitedNodes))
print(visitedNodes)
c) Breadth-first Search:
import queue as q
def breitensuche(adj,start,destination,visitedNodes=[]):
    queue=q.Queue()
    queue.put(start)
    while queue.qsize()>0:
        currentNode=queue.get()
        visitedNodes.append(currentNode)
        for successor in adj[currentNode]:
            if successor in visitedNodes:
                visitedNodes.append('Node'+str(successor)+
                    'was already visited earlier.')
                continue
            elif successor==destination:
                visitedNodes.append(successor)
                return True
            else:
                queue.put(successor)
    return False
adjList=f0:[1,2],1:[2,3],2:[4],3:[4,5],4:[],5:[] g
visitedNodes=[]
print(breitensuche(adjList,0,5,visitedNodes))
print(visitedNodes)

```

Spark

```

from operator import add
from pysparkling import Context
sc=Context()
a) Find 25 suppliers with the lowest account balance.
top25=(
    sc.textFile('supplier.tbl')
    .map(lambda line:line.split('j'))
    .map(lambda row:(row[1],float(row[5])))
    .sortBy(lambda row:row[1])
    .take(25)
)
# more efficient on PySpark (without global sorting, which causes substantial
# shuffle/repartitioning), exactly the same as above solution on Pysparkling,
# which implements top(num, key) as sortBy(key, ascending=False).take(num)
top25=(
    sc.textFile('supplier.tbl')
    .map(lambda line:line.split('j'))
    .map(lambda row:(row[1],float(row[5])))
    .top(25, key=lambda row:row[1])
    # or .takeOrdered(25, key=lambda row:row[1]) # only available on PySpark
)
print(top25)
b) How many suppliers have a positive account balance?
num pos bal=(
    sc.textFile('supplier.tbl')
    .map(lambda line:float(line.split('j')[5]))
    .filter(lambda acctbal:acctbal>0)
)

```

```
.count()
)
print(num pos bal)
```

c) Find out all brands produced by the same manufacturer and calculate the items number and the

total sales price for each brand of each manufacturer.

```
brandmfgcount=(
sc.textFile('part.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:((row[2],row[3]),1))
.reduceByKey(add)
)
brandmfgsum=(
sc.textFile('part.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:((row[2],row[3]),float(row[7])))
.reduceByKey(add)
)
1
result=brandmfgcount.join(brandmfgsum).sortBy(lambda x:x[0]).collect()
print(result)
#shortersolution
def addtupleselementwise(_args):
return tuple(sum(x)for xin zip(_args))
result=(
sc.textFile('part.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:((row[2],row[3]),(1,float(row[7]))))
.reduceByKey(addtupleselementwise)
)
print(result.sortBy(lambda x:x[0]).collect())
```

d) How many items have 3 words in their name?

```
num namelength3=(
sc.textFile('part.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:len(row[1].split()))
.filter(lambda length:length==3)
.count()
)
print(num namelength3)
```

e) How many different items does each supplier have?

```
supplier=(
sc.textFile('supplier.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:(row[0],(row[1],row[2],row[3],row[4],row[5],row[6])))
)
partsuppliers=(
sc.textFile('partsupp.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:(row[1],(row[0],row[2],row[3],row[4])))
)
supplierpsright=(
supplier.rightOuterJoin(partsuppliers)
.map(lambda x:(x[1][1][0],(x[1][0],x[0],x[1][1][1:])))
)
part=(
sc.textFile('part.tbl')
.map(lambda line:line.split('j'))
.map(lambda row:(row[0],(row[1],row[2],row[3],row[4],row[5],row[6],row[7],row[8])))
)
supplierpart=(
supplierpsright.rightOuterJoin(part)
.map(lambda x:(x[1][0][0][0],1))
```

```
.reduceByKey ( add )
)
print (supplierpart.collect())
2
```

Note: The is .join on Pysparkling, but it works incorrectly. The correct result for our query e) can be obtained using .rightOuterJoin, but in the general case that won't work.

It is possible to express and answer all 5 queries with Spark's RDDs. For the queries c) and e), however, the result is very elaborate and unattractive. DataFrame concept from Spark offers a better approach.

Exercise 5: PageRank with MapReduce (homework)

```
from mrsim import mr simulator
nodes = [ (('a', 0.2), ['a', 'c']),
          (('b', 0.2), ['a', 'd']),
          (('c', 0.2), ['b', 'c', 'd']),
          (('d', 0.2), ['c', 'e']),
          (('e', 0.2), []) ]
# this can be done in another map-reduce steps
# collect all vertex labels
# add all graph vertices to adjacency list of dead ends
def removeDeadends ( nodes ):
    vertexids = [ v[0][0] for v in nodes ]
    for v in nodes:
        adjacencylist = v[1]
        if not adjacencylist: # adjacencylist is empty
            adjacencylist.extend ( vertexids )
    return nodes
def map ( key, val ):
    res = []
    for link in val:
        res.append ( (link, ('val', key[1]/len (val))) )
        res.append ( (key[0], ('link', val)) )
    return res
def reduce ( key, val ):
    res = []
    prvalue = 0.0
    links = ()
    for v in val:
        if (v[0] == 'val'):
            prvalue = prvalue + v[1]
        else:
            links = v[1]
    prvalue = 0.8 * prvalue + (0.2)/len (nodes)
    res.append ( (key, prvalue), links )
    return res
intermediateResult = removeDeadends ( nodes )
# first 10 iterations
for i in range (1, 11):
    print ('Iteration' + str(i) + ':')
    intermediateResult = mr simulator (intermediateResult, map, reduce)
    print (intermediateResult)
```

Counting Triangles

```
import itertools
# tests, whether vertex x is less than vertex y (x < y) in graph;
# < defines strict total order on vertices
def vertexlt(x, y, graph):
    deg_x = len(graph[x]) # because all edges are bidirectional
    deg_y = len(graph[y])
    if deg_x < deg_y or deg_x == deg_y and x < y:
        return True
    return False
def counttriangles(graph):
    numtriangles = 0
    for v, adjacencylist in graph.items(): # vis vertex id
        # generator expression instead of list comprehension
        neighbors_of_greater_deg = (n for n in adjacencylist if vertexlt(v, n, graph))
        for u, w in itertools.combinations(neighbors_of_greater_deg, 2):
            if w in graph[u]: # because all edges are bidirectional
                print('Found triangle: ' + v + u + w)
                numtriangles += 1
    return numtriangles
graph = {'a': ['b', 'c', 'd', 'f'],
        'b': ['a', 'c', 'e'],
        'c': ['a', 'b', 'd', 'f', 'g'],
        'd': ['a', 'c', 'g'],
        'e': ['b', 'f'],
        'f': ['a', 'c', 'e'],
        'g': ['c', 'd']}
print(counttriangles(graph))
3
```

Pandas (Blatt 1 aufgabe 3)

a)

```
import pandas as pd
print('Find 25 Suppliers with lowest account balance.')

# csv laden und Spaltennamen vergeben:
supplier =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/supplier.tbl', sep='|',

names=['s_Id', 's_Name', 'Address', 'Nationkey', 'Phone', 'Acctbal', 's_Comment', 'Du
mmy'])

supplier.drop(supplier.columns[-1], axis=1, inplace=True)

top_25 = supplier.nsmallest(25, 'Acctbal')[['s_Name', 'Acctbal']]
# top_25 = (supplier.sort_values(by = 'Acctbal')[['s_Name', 'Acctbal']]).head
print(top_25)
```

b)

```

import pandas as pd
print('Wieviele Supps haben positive Balance?')

# csv laden und Spaltennamen vergeben:
supplier =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/supplier.tbl', sep='|',

names=['s_Id', 's_Name', 'Address', 'Nationkey', 'Phone', 'Acctbal', 's_Comment', 'Du
mmy'])

supplier.drop(supplier.columns[-1], axis=1, inplace=True)

supplier_pos_bal = supplier[supplier['Acctbal'] > 0.0][['s_Name']].count()

print(supplier_pos_bal)

supplier_pos_bal = supplier[supplier['Acctbal'] > 0.0]

print(supplier_pos_bal.shape[0])

```

c)

```

import pandas as pd
print('Find out all brands produced by the same manufacturer and calculate the
items number and thte total sales price for each')

# csv laden und Spaltennamen vergeben:
part =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/part.tbl', sep='|',

names=['p_Id', 'p_Name', 'Mfgr', 'Brand', 'Type', 'Size', 'Container', 'Retailprice',
'p_Comment', 'Dummy'])

part.drop(part.columns[-1], axis=1, inplace=True)

brand_manufacturer = part.groupby(['Mfgr', 'Brand']).agg({'p_Name':
'count', 'Retailprice': 'sum'})

print(brand_manufacturer)

```

d)

```

import numpy as np
import pandas as pd
print('How many items have 3 words in their name?')

```

```

# csv laden und Spaltennamen vergeben:
part =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/part.tbl', sep='|',

names=['p_Id', 'p_Name', 'Mfgr', 'Brand', 'Type', 'Size', 'Container', 'Retailprice',
'p_Comment', 'Dummy'])

part.drop(part.columns[-1], axis=1, inplace=True)

partNameSeries = part['p_Name']

print(partNameSeries[(partNameSeries.str.split()).map(len) ==
5].aggregate('count'))

#print(np.sum(partNameSeries.str.split().map(len) == 3))

```

e)

```

import pandas as pd
print('How many items have 3 words in their name?')

# csv laden und Spaltennamen vergeben:
supplier =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/supplier.tbl', sep='|',

names=['s_Id', 's_Name', 'Address', 'Nationkey', 'Phone', 'Acctbal', 's_Comment', 'Du
mmy'])

supplier.drop(supplier.columns[-1], axis=1, inplace=True)

part =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/part.tbl', sep='|',

names=['p_Id', 'p_Name', 'Mfgr', 'Brand', 'Type', 'Size', 'Container', 'Retailprice',
'p_Comment', 'Dummy'])

part.drop(part.columns[-1], axis=1, inplace=True)

partsupp =
pd.read_csv('C:/Dateien/Master/Skripte/2_Semester/Analyzing_Massive_Data_sets/
Exercises/Sheet01/partsupp.tbl', sep='|',

names=['p_FK', 's_FK', 'Availqty', 'Supplycost', 'Comment_partsupp', 'Dummy'])

partsupp.drop(supplier.columns[-1], axis=1, inplace=True)

```

```
supplierPS = pd.merge(supplier, partsupp, left_on = 's_Id', right_on = 's_FK')
supplierPart = pd.merge(supplierPS, part, left_on = 'p_FK', right_on = 'p_ID')

print(supplierPart.groupby('s_Name')['p_Name'].count())
```