

MAPPEDUCE CODE

Input Data

Each mapper processes a portion of the input data and each will be given a line

- Text is usually freeform so we use something like a regular expression

- sometimes is regular, like here

2012-01-01 12:01 San Jose 12.88 Amex

tab delimited, split by tab and extract values

... then use values needed to create key-value pairs for the ~~map~~ mapreduce job

Defining Mapper Code

```
def mapper():
```

```
    for line in sys.stdin:
```

```
        data = line.strip().split("\t")
```

```
        *date, time, store, item, cost, payment = data
```

we want to make sure that no matter what data we get, the mapper will keep working

```
        if len(data) == 6:
```

```
            ...
```

Between map and reduce there shuffle and sort!

Hadoop streaming allows you to write in ~~only~~ any language

Reducing

key-value pairs

~~Each~~ ~~map~~ ~~is~~ ~~split~~ among reducers.

Each will be sort (from ~~map~~ shuffle and sort)

Miami 12.34
Miami 99.07
Miami 3.14
Nyc 99.77
Nyc 88.98

```
def reducer:
```

```
    costTotal = 0
```

TO check when the store changes

```
    ddkey = None
```

```
    for line in sys.stdin:
```

```
        data = line.strip().split("\t")
```

wait for data

```
    if ddkey is set and  
    we finished processing  
    a store, print it  
    set variables for  
    next loop
```

```
    if len(data) != 2:  
        continue
```

```
    thiskey, thiscost = data
```

```
    if ddkey and ddkey != thiskey:  
        costTotal = 0  
        print ""
```

to print the result for the last store

```
    ddkey = thiskey  
    costTotal += float(thiscost)
```

```
    if ddkey != None:  
        print ""
```

testing it

```
read -50 purchases.txt > testFile  
cat testFile | ./mapper.py
```

for the reducer...

```
cat testFile | ./mapper.py | sort  
| ./reducer.py
```

```
hadoop -fs cat -cat  
output/part-0000.txt
```

↓
see the hadoop output!

you can run it on hadoop and check ~~program~~ program on the job tracker

MAPPING PATTERNS

- Filtering patterns, sampling, top-N
- summarization patterns, counting, min/max, statistics, index
- structural patterns, ~~combining~~ combining datasets

Filtering patterns

- Simple filter, ~~just~~ just a function
- Bloom Filter, ~~efficient~~ more efficient, probabilistic
- Sampling, note a smaller dataset from a larger by pulling out samples
- Random sampling
- top N

Summarization patterns

- ③ give quick easy high level understanding of your data
- Inverted Index
- Numerical summarizations (word/keyword count, mean, median, standard deviation)

Combiners (between map and reduce phase)

Transferring all the data from mappers to reducers is a lot of bandwidth, so we do some pre-reduction before passing the data to reducers

Structural Patterns

Migrating data to hadoop (take advantage of hierarchical data)

You need:

- Data sources linked by Foreign key
- Data must be structural and row-based

Don't need to do joins, save time reformatting your data by