

## Live Project 1

Retail Data Analysis using Pig and Hive

# Retail Transactions

- Retail stores daily generate millions of transactions logs.
- Analyzing these logs would generate beautiful insights and improve business.
- Storing these logs on traditional databases would be costly and scalability will be a big challenge.

# Volume of transactions

- Stores like walmart are spread across different locations.
- Daily millions of customers visit these stores and generate billions of logs.
- This billions of logs contribute to huge volume of data.

# Velocity of transactions

- In peak hours, 1000's of transactions will happen in any given second.
- 1000's of trasactions/sec across all stores contribute to high velocity.

# Variety of transactions

- The most widely known varieties of data generated by transactions:
  - Json Format
  - Xml Format
  - Csv Format

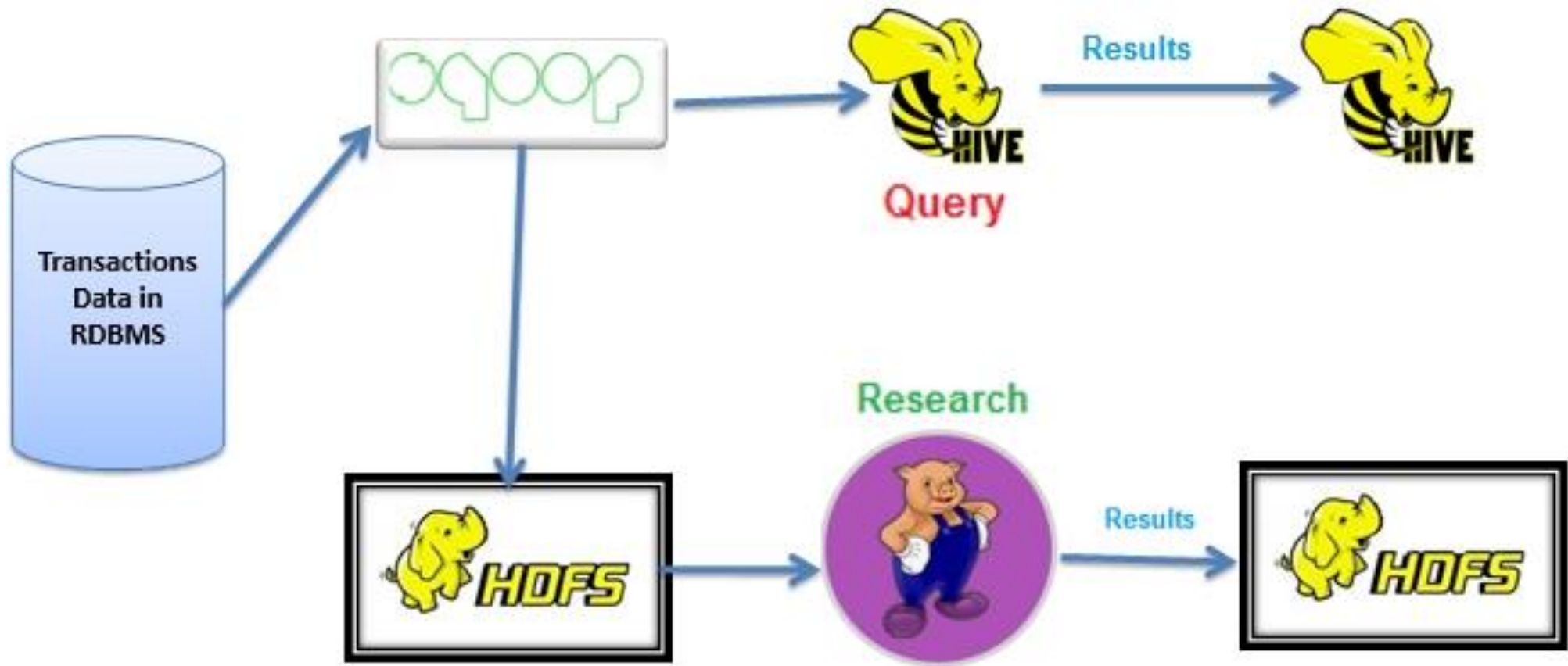
# Its Bigggg Data

- Having Huge Volume, High velocity and variety will make this data Big Data.
  - Challenges:
    - Storage
    - Scalability
    - Processing
    - Sharing

# Goal

- Process these logs over night as a batch to understand:
  - Demand of a given product.
  - Trend and seasonality of sales.
  - Understand performance of chain.
  - Loyal Customer identification.

# Solution





# Transactions dataset

Field	Type	Null	Key	Default	Extra
id	varchar(20)	YES		NULL	
chain	varchar(20)	YES		NULL	
dept	varchar(20)	YES		NULL	
category	varchar(20)	YES		NULL	
company	varchar(20)	YES		NULL	
brand	varchar(20)	YES		NULL	
date1	varchar(10)	YES		NULL	
productsize	int(11)	YES		NULL	
productmeasure	varchar(10)	YES		NULL	
purchasequantity	int(11)	YES		NULL	
purchaseamount	float	YES		NULL	

11 rows in set (0.00 sec)

# Load data into MySQL

```
-- select database
use training

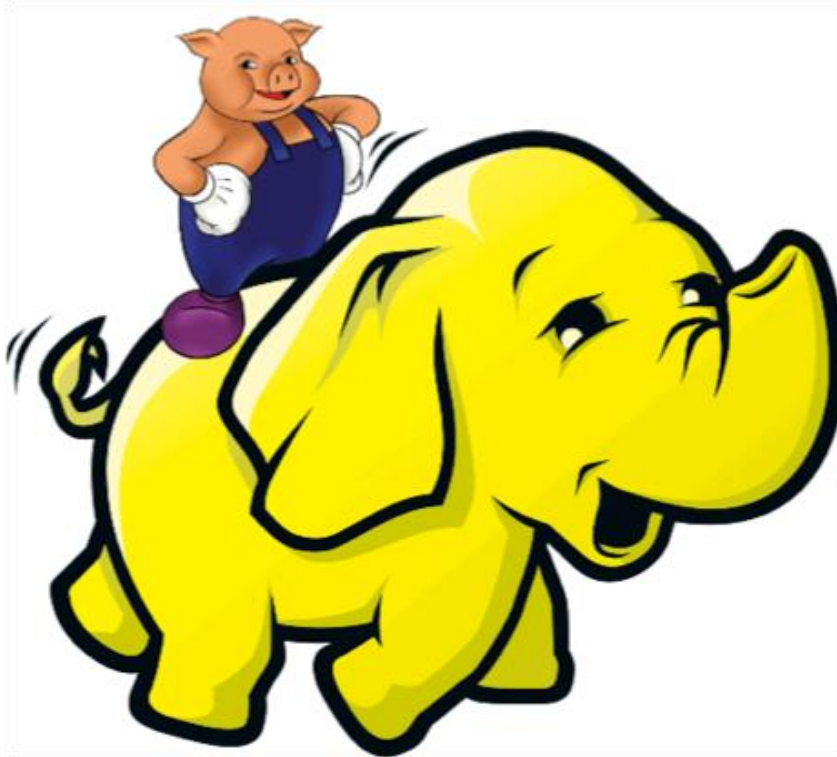
-- create table
CREATE TABLE transactions(id varchar(20),chain varchar(20), dept varchar(20),category varchar(20), company varchar(20),
brand varchar(20), date1 varchar(10), productsize int, productmeasure varchar(10), purchasequantity int, purchaseamount FLOAT);

-- load data
LOAD DATA INFILE '/home/cloudera/Desktop/transactions.csv' INTO TABLE transactions
FIELDS TERMINATED BY ',' ENCLOSED BY '"'
LINES TERMINATED BY '\r\n';

-- check

select * from transactions limit 5;
```

# Retail Data Analysis with Pig



- describe
- illustrate
- explain

**Use the above commands for better understanding of schema and data flow.**

# SQOOP data into HDFS

```
-- SQOOP data into HDFS

sqoop import --connect jdbc:mysql://localhost/training
--username training -P --table transactions -m 1 --target-dir TransactionsData
```

## Contents of directory [/user/training/TransactionsData](#)

Goto :

[Go to parent directory](#)

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
<a href="#">SUCCESS</a>	file	0 KB	1	64 MB	2015-04-05 21:52	rw-r--r--	training	supergroup
<a href="#">part-m-00000</a>	file	57.52 MB	1	64 MB	2015-04-05 21:52	rw-r--r--	training	supergroup

## Load data into pig

```
transactions = LOAD 'TransactionsData/part-m-00000' USING PigStorage(',') as  
(id:chararray,chain:chararray,dept:chararray,category:chararray,company:chararray,brand:ch  
ararray,date:chararray, productsize:float, productmeasure:chararray, purchasequantity:int,  
purchaseamount:float);
```

## Top 10 customers

```
custGroup = GROUP transactions BY id; --grouping
```

```
custSpending = FOREACH custGroup GENERATE group,  
SUM(transactions.purchaseamount) as spendings; --sum operation
```

```
custSpendingSort = ORDER custSpending BY spendings desc;
```

```
top10Cust = LIMIT custSpendingSort 10;
```

```
DUMP top10Cust;
```

```
STORE top10Cust INTO 'top10Cust'
```

## Snap shot of output

**File:** [/user/training/top10Cust/part-r-00000](#)

Goto :

[Go back to dir listing](#)

[Advanced view/download options](#)

```
86252    53592.899837329984
86246    52828.11981391907
57132131      20863.289817154408
83868868      15302.16985589452
67957375      14684.610046138987
58237989      14611.269816048443
91998805      14572.459874942899
94244413      14302.609862526879
54590006      13740.629864683375
56544981      13101.059877915308
```

```
chainGroup = GROUP transactions BY chain;
```

```
chainSales = FOREACH chainGroup GENERATE group,  
SUM(transactions.purchasequantity) as totalQuantity,  
SUM(transactions.purchaseamount) as totalSales;
```

```
dump chainSales;
```

```
STORE chainSales INTO 'chainSales';
```



# Snapshot of output

File: [/user/training/chainSales/part-r-00000](#)

Goto :

[Go back to dir listing](#)

[Advanced view/download options](#)

2	1194	3814.6399653851986
3	65062	208104.0280714687
4	245482	767159.3530194685
14	98427	365785.296934329
15	319263	1077677.4602135327
17	123889	394737.396487277
18	210325	693190.203932317
20	92454	291686.07754904777
58	4438	10223.759967654943
88	104758	324290.987072587
95	145003	441778.4061847329
205	34937	106421.01965124905

## Each chain, top 10 customers

```
chainGroupCust = GROUP transactions BY (chain,id);
```

```
chainGroupCustSpedings1 = FOREACH chainGroupCust GENERATE group, SUM(transactions.purchaseamount)  
as spendings;
```

```
chainGroupCustSpedings2= FOREACH chainGroupCustSpedings1 generate group.chain as  
chain,group.id as id, spendings;
```

```
chainGroupCustSpedings3= GROUP chainGroupCustSpedings2 BY chain;
```

```
chainTop10Cust = FOREACH chainGroupCustSpedings3{  
chainGroupCustSpedingsSort = ORDER chainGroupCustSpedings2 BY spendings DESC;  
top10Cust = LIMIT chainGroupCustSpedingsSort 10;  
GENERATE top10Cust;  
}
```

```
chainTop10Cust = FOREACH chainTop10Cust GENERATE FLATTEN(top10Cust);
```

```
STORE chainTop10Cust INTO 'chainTop10Cust '
```

# Snapshot of output

File: [/user/training/chainTop10Cust/part-r-00000](#)

Goto :

[Go back to dir listing](#)

[Advanced view/download options](#)

14	67957375	14684.610046138987
14	58237989	14611.269816048443
14	17456470	12878.939868450165
14	13074629	12329.949950814247
14	20552767	11653.639915667474
14	97110346	11154.869882548228
14	90548180	10474.799921853468
14	55304581	10256.949917048216
14	59816528	9873.22992990911
14	91883639	8436.549869716167
15	57132131	20863.289817154408
15	54590006	13740.629864683375
15	93752475	12884.169838543981
15	51531106	12739.759873157367
15	13251776	12512.589874466881
15	85851216	12289.38995597884

**This operation cannot be done using hive queries**

## Each customer most bought brand

```
CustBrandGroup = GROUP transactions BY (id,brand);
```

```
CustBrandQuantity = FOREACH CustBrandGroup GENERATE group,  
SUM(transactions.purchasequantity) as sales;
```

```
CustBrandQuantity = FOREACH CustBrandQuantity GENERATE  
group.brand as brand, group.id as id, sales;
```

```
CustBrandQuantityGroup = GROUP CustBrandQuantity BY brand;
```

```
custTop5Brands = FOREACH CustBrandQuantityGroup{  
    CustBrandQuantityGroupSort = ORDER CustBrandQuantity BY sales DESC;  
    top5Brand = LIMIT CustBrandQuantityGroupSort 5;  
    GENERATE top5Brand;}
```

```
custTop5Brands = FOREACH custTop5Brands GENERATE FLATTEN(top5Brand);
```

```
STORE custTop5Brands INTO 'custTopFiveBrands';
```

# Snapshot of output

**File:** [/user/training/custTopFiveBrands/part-r-00000](#)

Goto :

[Go back to dir listing](#)  
[Advanced view/download options](#)

[View Next chunk](#)

101609	01410024	1
101609	22326006	1
101616	62407885	7
101616	20552767	5
101616	97110346	3
101616	58237989	2
101616	59816528	2
101646	84841861	1
101646	81242696	1
10165	49806426	79
10165	70401965	59
10165	76895949	56
10165	95508708	51
10165	86252	48

**This operation cannot be done using hive queries**

# Top 10 brands

```
brandGroup = GROUP transactions BY brand;
```

```
brandPurchase = FOREACH brandGroup GENERATE group,  
SUM(transactions.purchaseamount) as purchase;
```

```
brandPurchaseSort = ORDER brandPurchase BY purchase desc;
```

```
top10Brands = LIMIT brandPurchaseSort 10;
```

```
STORE top10Brands INTO 'top10Brands';
```

# Snapshot of output

**File:** [/user/training/top10Brands/part-r-00000](#)

Goto :

[Go back to dir listing](#)

[Advanced view/download options](#)

[View Next chunk](#)

```
15704 449154.8166325502
10786 221273.3691299595
0      102354.53942238353
12908 61599.71946503781
16397 52987.138719622046
30626 44524.759959416464
13310 37686.129275966436
33170 36812.93964109942
17286 36278.89945333451
```

## top 10 companies

```
companyGroup = GROUP transactions BY company;
```

```
companyPurchase = FOREACH companyGroup GENERATE group,  
SUM(transactions.purchaseamount) as purchase;
```

```
companyPurchaseSort = ORDER companyPurchase BY purchase desc;
```

```
top10Companies = LIMIT companyPurchaseSort 10;
```

```
STORE top10Companies INTO 'top10Companies';
```



# Snapshot of output

File: [/user/training/top10Companies/part-r-00000](#)

Goto :

[Go back to dir listing](#)

[Advanced view/download options](#)

102113020	761962.2238900699
107989373	136652.5689506065
103700030	115869.42819562927
107675070	109698.44996777177
102840020	76357.61965951882
104900040	58961.50925568119
101600010	58111.499768570065
104400040	51936.17975332588
101200010	51542.929359253496
103800030	47719.35976731032

## Chain Year Monthly Sales

```
chainYearMonSales = FOREACH transactions GENERATE  
chain,STRSPLIT(date,'/',3),purchaseamount as sales;
```

```
chainYearMonSales = FOREACH chainYearMonSales GENERATE chain, $1.$0 as  
month, $1.$2 as year, sales;
```

```
chainYearMonSalesGroup = GROUP chainYearMonSales by (chain,year,month);
```

```
chainYearMonGroupSales = FOREACH chainYearMonSalesGroup GENERATE group,  
SUM(chainYearMonSales.sales) as totalsales;
```

```
chainYearMonGroupSales = FOREACH chainYearMonGroupSales GENERATE  
group.chain as chain, group.year as year, group.month as month, totalsales;
```

```
STORE chainYearMonGroupSales INTO 'chainYearMonGroupSales';
```

# Snapshot of output

File: <a href="#">/user/training/chainYearMonGroupSales/part-r-00000</a>			
Goto : <input type="text" value="/user/training/chainYearMor"/> <input type="button" value="go"/>			
<a href="#">Go back to dir listing</a>			
<a href="#">Advanced view/download options</a>			
2	2012	3	81.23999857902527
2	2012	4	208.37999892234802
2	2012	5	62.70999884605408
2	2012	6	176.60999804735184
2	2012	7	208.92999720573425
2	2012	8	308.0299973487854
2	2012	9	217.59999758005142
2	2012	10	386.1499967575073
2	2012	11	491.4499979019165
2	2012	12	481.7199950516224
2	2013	1	685.2599949836731
2	2013	2	139.96999841928482
2	2013	3	366.5899957418442
3	2012	3	13554.699863106012
3	2012	4	11968.10988478735
3	2012	5	13221.750001200363



## Sqoop data into hive

```
-- sqoop data from mysql to hive

sqoop import --connect jdbc:mysql://localhost/training \
  --username training -P \
  --table transactions \
  --hive-import \
  --hive-table transactions_staging -m 1
```

# Partitioning and bucketing

- Partitioning and bucketing in hive will let you do faster querying.
- For dynamic partitioning, load the data in to staging table which is already done.
- Now create a production table, and insert data.

## Open hive and set below properties

```
-- Partitioning and bucketing  
set hive.enforce.bucketing = true;  
SET hive.exec.dynamic.partition = true;  
SET hive.exec.dynamic.partition.mode = nonstrict;
```

## Create production table and load data

```
-- production table
CREATE TABLE transactions_production
( id string,
  dept string,
  category string,
  company string,
  brand string,
  date1 string,
  productsize int,
  productmeasure string,
  purchasequantity int,
  purchaseamount double)
PARTITIONED BY (chain string) CLUSTERED BY(id) INTO 5 BUCKETS
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE;

INSERT OVERWRITE TABLE transactions_production PARTITION (chain)
select id,dept, category, company,brand,date1,productsize,productmeasure,
purchasequantity,purchaseamount,chain from transactions_staging;
```



# Partitioning snapshot

Goto :

[Go to parent directory](#)

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
<a href="#">chain=14</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=15</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=17</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=18</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=2</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=20</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=205</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=3</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=4</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=58</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=88</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive
<a href="#">chain=95</a>	dir				2015-04-04 11:08	rw-rw-rw-	cloudera	hive

# Bucketing snapshot in chain=14

Contents of directory [/user/hive/warehouse/transactions\\_production/chain=14](#)

Goto :

[Go to parent directory](#)

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
<a href="#">000000_0</a>	file	820.55 KB	3	128 MB	2015-04-04 11:07	rw-rw-rw-	cloudera	hive
<a href="#">000001_0</a>	file	873.68 KB	3	128 MB	2015-04-04 11:07	rw-rw-rw-	cloudera	hive
<a href="#">000002_0</a>	file	632.47 KB	3	128 MB	2015-04-04 11:07	rw-rw-rw-	cloudera	hive
<a href="#">000003_0</a>	file	675.31 KB	3	128 MB	2015-04-04 11:07	rw-rw-rw-	cloudera	hive
<a href="#">000004_0</a>	file	942.64 KB	3	128 MB	2015-04-04 11:08	rw-rw-rw-	cloudera	hive

## Top 10 customers

```
-- Top 10 customers

select id, sum(purchaseamount) as custSpendings from transactions_production
group by id
sort by custSpendings DESC
limit 10;
```

## Snapshot of output

```
86252      53592.899999999957
86246      52828.119999999952
57132131    20863.290000000036
83868868    15302.1699999999785
67957375    14684.609999999933
58237989    14611.2699999999817
91998805    14572.459999999987
94244413    14302.609999999986
54590006    13740.6299999999846
56544981    13101.0599999999881
Time taken: 178.254 seconds
```

## Chain wise sales

```
-- chain wise sales  
select chain, sum(purchaseamount), sum(purchasequantity) from transactions_production  
group by chain;
```

## Snapshot of output

```
OK
14      365785.2999998729      98427
15      1077677.4699990302     319263
17      394737.39999981265     123889
18      693190.209999508       210325
2       3814.6399999999644      1194
20      291686.0799999232      92454
205     106421.01999999923      34937
3       208104.02999998012      65062
4       767159.3599993604      245482
58      10223.759999999922      4438
88      324290.9899998967      104758
95      441778.4099997955      145003
Time taken: 62.781 seconds
```

## Top 10 brands

```
-- top 10 brands
select brand, sum(purchaseamount) as custSpending from transactions_production
group by brand
sort by custSpending DESC
limit 10;
```

## Snapshot of output

```
OK
15704      449154.81999997314
10786      221273.369999998149
0          102354.540000000849
12908      61599.719999999789
16397      52987.139999999947
30626      44524.760000000001
13310      37686.1300000001845
33170      36812.9400000002265
17286      36278.9000000000904
9886       30633.2600000002185
Time taken: 141.856 seconds
```



## Top 10 companies

```
-- top 10 companies  
select company, sum(purchaseamount) as custSpendings from transactions_production  
group by company  
sort by custSpendings DESC  
limit 10;
```

## Snapshot of output

```
Total MapReduce CPU Time Spent: 19 seconds 760 msec
OK
102113020      761962.2299992407
107989373      136652.57000001118
103700030      115869.43000001504
107675070      109698.45000000006
102840020      76357.6199999992
104900040      58961.5100000003
101600010      58111.500000000546
104400040      51936.17999999996
101200010      51542.92999999974
103800030      47719.360000000524
Time taken: 134.221 seconds
```

## Chain Year Monthly Sales

```
-- Chain Year Monthly Sales  
select chain, split(date1,'/')[2] as year1, split(date1,'/')[0] as month1,  
sum(purchaseamount) as totalsales from transactions_production  
group by chain,split(date1,'/')[0],split(date1,'/')[2];
```

# Snapshot of output

```
OK
14      2013      1      25989.1000000000908
14      2012     10      24914.2200000000976
14      2012     11      26941.650000000088
14      2012     12      34061.44000000014
14      2013      2      24553.6600000000826
14      2012      3      22811.4000000000634
14      2013      3      25149.640000000094
14      2012      4      23113.8900000000614
14      2013      4      17588.639999999994
14      2012      5      26469.710000000088
14      2013      5      10979.3399999999844
14      2012      6      24692.900000000072
14      2013      6      8059.4299999999861
14      2012      7      21334.340000000049
14      2013      7      543.80000000000003
14      2012      8      23286.670000000069
14      2012      9      25295.470000000095
15      2013      1      77346.020000000096
15      2012     10      77011.470000000112
15      2012     11      76465.900000000053
15      2012     12      97227.450000000458
15      2013      2      70020.989999999887
15      2012      3      69659.749999999933
15      2013      3      77026.0800000001
```

# Storing results

- To store results

```
-- store results

CREATE TABLE target
AS
SELECT col1,col2
FROM SOURCE
```

```
-- Example create table top10companies
CREATE TABLE top10companies
AS
select company, sum(purchaseamount) as custSpending from transactions_production
group by company
sort by custSpending DESC
```