

# RBDA Project

## Project - Rachit Jain - rj2219

---

### Explanation

---

#### Dataset description:

1. The dataset that I used is NY citywide payroll dataset for a fiscal year [1]. This data was collected to study how the state's budget was being allocated.
2. The data has data points from fiscal year 2014 to 2020.
3. The data consists of the following columns:

Column Name	Description	Preprocessed	Kept/Dropped/Converted
Fiscal Year	Fiscal year for the data point	No	Kept
Payroll number	Agency number	No	Dropped
Agency Name	Name of the agency which was making payments	The agency names had commas which needed to be preprocessed so that they could be read in the <a href="#">mapper</a> .	Kept
Last Name	Last name of the employee	No	Dropped
First Name	First name of the employee	No	Dropped
Mid Init	Middle initial of the employee	No	Dropped
Agency Start Date	Date on which employee began working	No	Dropped

Column Name	Description	Preprocessed	Kept/Dropped/Converted
Work Location Borough	Borough of employee's work location	Removed boroughs other than NYC boroughs - had boroughs like Washington DC, Albany, etc.	Kept
Title Description	Civil service title description of the employee	The agency names had commas which needed to be preprocessed so that they could be read in the <code>mapper</code> .	Dropped
Leave Status as of June 30th	Status of employment	No	Dropped
Base Salary	Base salary of the employee	Some salaries were not in <code>per annum</code> calculation. Converted those.	Converted
Pay Basis	When the employee is paid	Used this in calculation for getting total salary.	Converted
Regular Hours	Number of hours worked	No. Data issue	Dropped
Regular Gross Paid	Amount paid to the employee	No. Data issues	Dropped
OT Hours	Number of overtime hours worked	No	Kept
Total OT Paid	Total overtime paid	No	Converted
Total Other Pay	Salary in addition to base	No	Converted

4. The data set had a few issues where some columns did not have the correct data.
5. Total rows:
  - a. Before profiling and cleaning: 3923291

- b. After profiling and cleaning: 2470563

### **Data Profiling and Cleaning:**

1. The code has been uploaded to [GitHub](#) in a private repository [2].
2. The overall idea for the data is to drop the columns that are not required and profile the data to get the overall `totalPay` of a person which we will later relate to the purchasing capacity of the borough.
3. I set the `mapred.textoutputformat.separator` parameter as `,` as I need to write the data in `.csv` format.
4. My input data is in the `.csv` format and output data is being written as `key -> processed_row`.
5. In my case, the `key` corresponds to the borough since all our analysis is borough based and the `processed_row` corresponds to the row of the csv after the columns are dropped and the row has been profiled and cleaned.
6. The mapper has a `set` of boroughs which are used for the processing -
 

```
"MANHATTAN", "BRONX", "BROOKLYN", "QUEENS", "STATEN ISLAND" .
```
7. The first step of cleaning the data was reading the input file which had `,` in the descriptions. This was solved by checking if the text was in a `string` format indicating one column in the data and the comma was converted to a semi colon. (`,` → `;`).
8. The above step made the parsing of the row based on a comma and we could easily divide the arguments.
9. If the borough was not an NYC borough, it was not included in the final HDFS cluster and the row was dropped.
10. I checked the `payBasis` and based on that, I converted each calculation into `per annum`.
11. Based on heuristics in [3], the following calculations were made:
  - a. `per annum or prorated annual` → `base salary + total OT paid + other allowances`
  - b. `per day` → `(base salary * number of days * working weeks) + total OT paid + other allowances`
  - c. `per hour` → `(base salary * number of hours * number of days * working weeks) + total OT paid + other allowances`

12. The average working hours are 7 and the average number of working weeks are 49.
13. Based on the above, the total pay of the person was calculated.
14. In the reducer, the output was going to be of the same format as that from the mapper but the only difference would be that I am dropping the rows which have a pay of less than 20000.
15. Minimum wage in NYC is \$14.20 per hour by end of December 2022 [4].
16. Using this for the calculation:
  - a.  $14.20 * 7 * 5 * 49 = 24,453$
  - b. But the minimum wage has been lesser in the previous years.
  - c. Using this, I came up with the heuristic of removing the total salaries of < 20,000 since that data might be incorrectly recorded.
17. The rows that were dropped in the reducer were used to form a statistic as to how many data points could be incorrect during the parsing of the data.
18. The dropped rows were aggregated and the final count of the borough\_year were written to the HDFS.
19. This could give us more insights as to how the data was recorded and how many rows could finally be used for the processing.
20. The final outputs were written to the HDFS with the type Text, Text.

## Screenshots

---

```

● ○ ● rachitjain@Rachits-MacBook-Pro:~/Desktop/Courses/Sem3/RBDA... ✘ 2
Last login: Wed Nov 16 18:02:25 on ttys001
(base) ➔ Project gcloud compute scp --recurse /Users/rachitjain/Desktop/Courses/Sem3/RBDA/Assignments/Project/citywide-payroll-data-fiscal-year.csv rj2219_nyu_edu@nyu-dataproc-m:lab3/ --project hpc-dataproc-19b8 --zone us-central1-f
External IP address was not found; defaulting to using IAP tunneling.
WARNING:

To increase the performance of the tunnel, consider installing NumPy. For instructions,
please see https://cloud.google.com/iap/docs/using-tcp-forwarding#increasing\_the\_tcp\_upload\_bandwidth

citywide-payroll-data-fiscal-year.csv          100%  583MB   2.2MB/s   04:19

Updates are available for some Google Cloud CLI components. To install them,
please run:
$ gcloud components update

(base) ➔ Project

```

`scp` command for pushing data to the dataproc cluster

```

● ○ ● rj2219_nyu_edu@nyu-dataproc-m: ~/project ✘ 1
rj2219_nyu_edu@nyu-dataproc-m:~/project$ hadoop fs -put citywide-payroll-data-fiscal-year.csv project/
rj2219_nyu_edu@nyu-dataproc-m:~/project$ hadoop fs -ls project/
Found 1 items
-rw-r--r--  1 rj2219_nyu_edu rj2219_nyu_edu  611400234 2022-11-16 23:11 project/citywide-payroll-data-fiscal-year.csv
rj2219_nyu_edu@nyu-dataproc-m:~/project$ 

```

Showing that the dataset has been pushed onto the `hadoop` file system

```
Last login: Sat Nov 26 11:00:12 on ttys002
(base) + ~ gcloud compute scp --recurse /Users/rachitjain/Desktop/Courses/Sem3/RBDA/Assignments/Project/nyc-crime-arrest-payroll rj2219_nyu_edu@nyu-dataproc-m:project/ --project hpc-dataproc-19b8 --zone us-central1-f
External IP address was not found; defaulting to using IAP tunneling.
WARNING:
To increase the performance of the tunnel, consider installing NumPy. For instructions,
please see https://cloud.google.com/iap/docs/using-tcp-forwarding#increasing_the_tcp_upload_bandwidth

PayrollFiscal.java          100% 1292   27.7KB/s  00:00
Project_Abstract.pdf        100% 64KB  449.2KB/s  00:00
README.md                    100% 26    0.8KB/s  00:00
PayrollFiscalReducer.java   100% 1516   36.8KB/s  00:00
PayrollFiscalMapper.java    100% 2030   49.9KB/s  00:00
.githooks                   100% 278    6.6KB/s  00:00
config                      100% 328    7.5KB/s  00:00
b502183ce3aa45919c3461e7f8a8bb50c1e375c 100% 590   13.7KB/s  00:00
5a356ca043f13fe93f391a168bb8ef6169c53a 100% 504   12.1KB/s  00:00
f1b86687f30ff91e58a1264dd8131f03ef0a62 100% 236   5.6KB/s  00:00
pack-bc98f871d55994e1ff37b058259b1211bf53a3c 100% 1184   26.2KB/s  00:00
pack-bc98f871d55994e1ff37b058259b1211bf53a3c 100% 868   20.9KB/s  00:00
49e4e65b35087e5573701766122c6973d40e 100% 58KB  469.0KB/s  00:00
06702ed62a4ca09fd973f36ca342e284d7714c 100% 170    4.0KB/s  00:00
fe990e3eaacbf834dd1b74b79de08.d24fe12 100% 752   17.9KB/s  00:00
HEAD                         100% 21    0.5KB/s  00:00
exclude                     100% 240   5.8KB/s  00:00
HEAD                         100% 379   9.4KB/s  00:00
main                         100% 379   9.1KB/s  00:00
HEAD                         100% 214   5.1KB/s  00:00
main                         100% 154   3.5KB/s  00:00
description                  100% 73    1.9KB/s  00:00
commit-msg.sample            100% 896   21.7KB/s  00:00
pre-rebase.sample            100% 4898  101.9KB/s  00:00
pre-commit.sample            100% 1643   35.8KB/s  00:00
applypatch-msg.sample       100% 478   11.3KB/s  00:00
fsmonitor-watchman.sample   100% 4655  99.9KB/s  00:00
pre-receive.sample           100% 544   13.2KB/s  00:00
prepare-commit-msg.sample   100% 1492  34.6KB/s  00:00
post-update.sample           100% 189   4.4KB/s  00:00
pre-merge-commit.sample    100% 416   9.7KB/s  00:00
pre-applypatch.sample      100% 424   9.5KB/s  00:00
pre-push.sample              100% 1374  30.3KB/s  00:00
update.sample                100% 3650  85.3KB/s  00:00
push-to-checkout.sample     100% 2783  63.2KB/s  00:00
main                         100% 41    1.0KB/s  00:00
HEAD                         100% 50    0.9KB/s  00:00
main                         100% 41    1.2KB/s  00:00
index                        100% 569   13.8KB/s  00:00
packed-refs                  100% 112   2.7KB/s  00:00
COMMIT_EDITMSG
```

Code files being pushed to the cluster. The extra files here are the files due to `git`.

```
● ● ● rj2219_nyu_edu@nyu-dataproc-m: ~/project/nyc-crime-arrest-pa... ↵
rj2219_nyu_edu@nyu-dataproc-m:~/project/nyc-crime-arrest-payroll$ javac -classpath `hadoop classpath` PayrollFiscalMapper.java
rj2219_nyu_edu@nyu-dataproc-m:~/project/nyc-crime-arrest-payroll$ javac -classpath `hadoop classpath` PayrollFiscalReducer.java
rj2219_nyu_edu@nyu-dataproc-m:~/project/nyc-crime-arrest-payroll$ javac -classpath `hadoop classpath`:. PayrollFiscal.java
rj2219_nyu_edu@nyu-dataproc-m:~/project/nyc-crime-arrest-payroll$ jar cvf PayrollFiscal..jar *.class
added manifest
adding: PayrollFiscal.class(in = 1654) (out= 906)(deflated 45%)
adding: PayrollFiscalMapper.class(in = 2664) (out= 1239)(deflated 53%)
adding: PayrollFiscalReducer.class(in = 2592) (out= 1150)(deflated 55%)
rj2219_nyu_edu@nyu-dataproc-m:~/project/nyc-crime-arrest-payroll$
```

Setting `classpath` and creating the `.jar` on hadoop.

```
rj2219_nyu_edu@nyu-dataproc-m:~/project/nyc-crime-arrest-payroll$ hadoop jar PayrollFiscal.jar PayrollFiscal project/citywide-payroll-data-fiscal-year.csv project/output
2022-11-27 19:37:53.010 INFO client.RMProxy: Connecting to ResourceManager at nyu-dataproc-m/192.168.1.39:8032
2022-11-27 19:37:53.181 INFO client.ANSProxy: Connecting to Application History server at nyu-dataproc-m/192.168.1.39:10200
2022-11-27 19:37:53.485 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2022-11-27 19:37:53.502 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/rj2219_nyu_edu/.staging/job_1668611926937_1075
2022-11-27 19:37:53.748 INFO input.FileInputFormat: Total input files to process : 1
2022-11-27 19:37:53.945 INFO Configuration.deprecation: mapred.textoutputformat.separator is deprecated. Instead, use mapreduce.output.textoutputformat.separator
2022-11-27 19:37:54.043 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1668611926937_1075
2022-11-27 19:37:54.044 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-11-27 19:37:54.223 INFO conf.Configuration: resource-types.xml not found
2022-11-27 19:37:54.223 INFO resource.ResourceUtils: Unable to find "resource-types.xml".
2022-11-27 19:37:54.390 INFO mapreduce.Job: Submitted application application_1668611926937_1075
2022-11-27 19:37:54.390 INFO mapreduce.Job: Running job: job_1668611926937_1075
2022-11-27 19:38:02.546 INFO mapreduce.Job: Job job_1668611926937_1075 running in uber mode : false
2022-11-27 19:38:02.548 INFO mapreduce.Job: map 0% reduce 0%
2022-11-27 19:38:11.657 INFO mapreduce.Job: map 20% reduce 0%
2022-11-27 19:38:12.663 INFO mapreduce.Job: map 40% reduce 0%
2022-11-27 19:38:13.668 INFO mapreduce.Job: map 100% reduce 0%
2022-11-27 19:38:27.738 INFO mapreduce.Job: map 100% reduce 100%
2022-11-27 19:38:28.750 INFO mapreduce.Job: Job job_1668611926937_1075 completed successfully
2022-11-27 19:38:28.838 INFO mapreduce.Job: Counters
File System Counters
  FILE: Number of bytes read=204221036
  FILE: Number of bytes written=409919379
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=611417378
  HDFS: Number of bytes written=147370588
  HDFS: Number of read operations=20
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=3
  HDFS: Number of bytes read erasure-coded=0
Job Counters
  Killed map tasks=1
  Launched map tasks=5
  Launched reduce tasks=1
  Rack-local map tasks=5
  Total time spent by all maps in occupied slots (ms)=171216
  Total time spent by all reduces in occupied slots (ms)=43980
  Total time spent by all map tasks (ms)=42804
  Total time spent by all reduce tasks (ms)=1095
  Total vcore-milliseconds taken by all map tasks=42804
  Total vcore-milliseconds taken by all reduce tasks=1095
  Total megabyte-milliseconds taken by all map tasks=175325184
  Total megabyte-milliseconds taken by all reduce tasks=45035520
Map-Reduce Framework
```

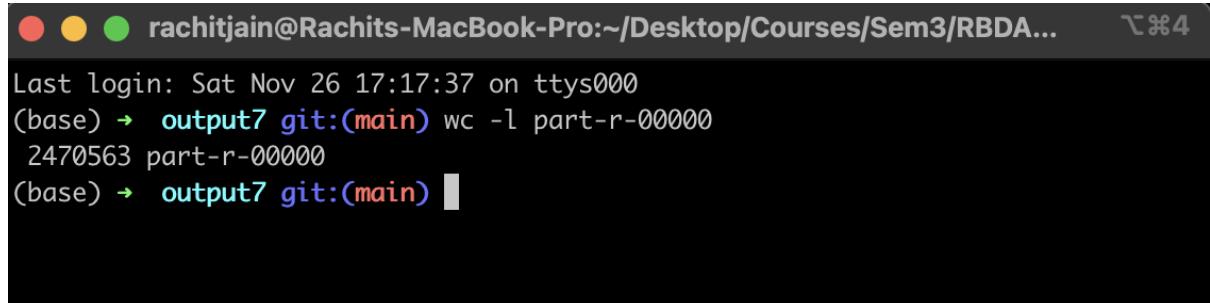
Screenshot showing running of the `mapreduce` task.

```
(base) → nyc-crime-arrest-payroll git:(main) ✘ git status
On branch main
Your branch is up to date with 'origin/main'.

Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git restore <file>..." to discard changes in working directory)
    modified:   .gitignore
    modified:   PayrollFiscal.java
    modified:   PayrollFiscalMapper.java
    modified:   PayrollFiscalReducer.java

no changes added to commit (use "git add" and/or "git commit -a")
(base) → nyc-crime-arrest-payroll git:(main) ✘ git add .
(base) → nyc-crime-arrest-payroll git:(main) ✘ git commit -m "Implemented Payro
llFiscal"
[main 434849b] Implemented PayrollFiscal
 4 files changed, 25 insertions(+), 22 deletions(-)
(base) → nyc-crime-arrest-payroll git:(main) git push origin main
Enumerating objects: 11, done.
Counting objects: 100% (11/11), done.
Delta compression using up to 8 threads
Compressing objects: 100% (6/6), done.
Writing objects: 100% (6/6), 1001 bytes | 1001.00 KiB/s, done.
Total 6 (delta 4), reused 0 (delta 0), pack-reused 0
remote: Resolving deltas: 100% (4/4), completed with 4 local objects.
To https://github.com/rachitjain2706/nyc-crime-arrest-payroll.git
 f6067b2..434849b main -> main
(base) → nyc-crime-arrest-payroll git:(main) █
```

`git status` output shown and data pushed to git.

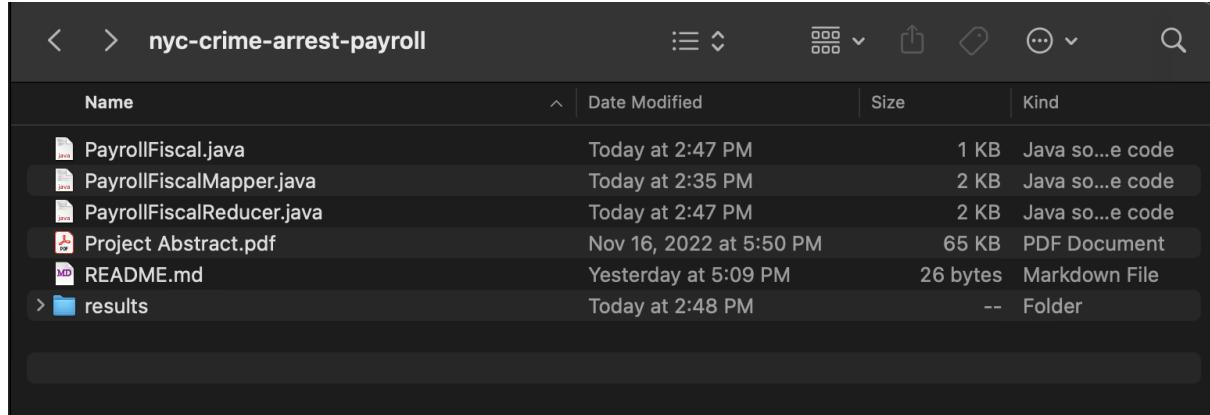


```

rachitjain@Rachits-MacBook-Pro:~/Desktop/Courses/Sem3/RBDA... ⌘⌘4
Last login: Sat Nov 26 17:17:37 on ttys000
(base) ➔ output7 git:(main) wc -l part-r-00000
2470563 part-r-00000
(base) ➔ output7 git:(main)

```

Output file after the mapreduce task ends. Shows the number of records kept after preprocessing.



Directory structure

A1	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
354410	2017	COMMUNITY COLLEGE (HOSTOS HIGSMITH WILLIAM				2008-08-11T BRONX		CUNY CUSTC CEASED	31329 per Annum	0	1737.25	0	66.77		17.95	
354411	2017	COMMUNITY COLLEGE (HOSTOS HICKMAN TIFFANY M				2015-07-01T BRONX		ADJUNCT LEI SEASONAL	194.52 per Day	0	258.72	0	0	0	0	
354412	2017	COMMUNITY COLLEGE (HOSTOS HICKMAN TIFFANY M				2015-07-01T BRONX		CONTINUING ACTIVE	36.64 per Hour	0	156.91	0	0	0	0	
354413	2017	COMMUNITY COLLEGE (HOSTOS HIDALGO PA CARLOS M				2005-11-15T BRONX		COLLEGE ASI ACTIVE	15.53 per Hour	0	225.9	0	0	0	1.16	
354414	2017	COMMUNITY COLLEGE (HOSTOS HILOAD LEE JACOB B				2013-02-06T BRONX		COLLEGE ASI CEASED	10.58 per Hour	0	369.9	0	0	0	133.41	
354415	2017	COMMUNITY COLLEGE (HOSTOS HERRERA YUNIOR				2010-10-25T BRONX		CONTINUING ACTIVE	36.64 per Hour	120	5627.98	0	0	0	0	
354416	2017	COMMUNITY COLLEGE (HOSTOS HERRERA DE KAILYN A				2017-02-10T BRONX		COLLEGE ASI ACTIVE	12 per Hour	307.75	3279	0	0	0	0	
354417	2017	COMMUNITY COLLEGE (HOSTOS HERRERA WILMER				2015-09-08T BRONX		COLLEGE ASI ACTIVE	12 per Hour	1001.25	12115.06	0	0	0	725.6	
354418	2017	COMMUNITY COLLEGE (HOSTOS HIAMANG TOKUNBO M				2013-09-23T BRONX		ADJUNCT LEI SEASONAL	59.7 per Day	80	6201.81	0	0	0	0	
354419	2017	COMMUNITY COLLEGE (HOSTOS HIGANDE EVELYN N				2017-09-06T BRONX		ASSISTANT TA ACTIVE	45957 per Annum	730	18371.03	0	0	0	0	
354420	2017	COMMUNITY COLLEGE (HOSTOS HOLMES ELAINE L				1995-09-22T BRONX		CUNY OFFICE CEASED	38407 per Annum	890.25	29721.91	0	0	0	1033	
354421	2017	COMMUNITY COLLEGE (HOSTOS HOLMES NASHANA C				2012-10-27T BRONX		CONTINUING ACTIVE	36.64 per Hour	0	45.55	0	0	0	0	
354422	2017	COMMUNITY COLLEGE (HOSTOS IFESANWA ADAMSON A				2011-06-29T BRONX		ADJUNCT LEI SEASONAL	74.32 per Day	160	11261.99	0	0	0	0	
354423	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ NORBERTO				2010-01-04T BRONX		ADJUNCT LEI SEASONAL	126.41 per Day	0	3017.31	0	0	0	1000	
354424	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ NORBERTO				2010-01-04T BRONX		COLLEGE ASI ACTIVE	18.27 per Hour	142.75	2493.35	0	0	0	22.85	
354425	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ LUIS				2015-01-26T BRONX		ADJUNCT AS SEASONAL	55.15 per Day	176.43	17259.75	0	0	0	328	
354426	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ NORBERTO				2013-01-07T BRONX		NON-TEACH ACTIVE	44.66 per Hour	0	2637.08	0	0	0	0	
354427	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ TERESA				2016-09-12T BRONX		CUNY OFFICE CEASED	29497 per Annum	455	7954.19	0	0	0	44.75	
354428	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ TERESA				2016-12-23T BRONX		COLLEGE ASI ACTIVE	14 per Hour	724	9411.5	0	0	0	17.5	
354429	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ AMARIS F				2012-07-02T BRONX		COLLEGE ASI ACTIVE	12.13 per Hour	0	227.47	0	0	0	13.55	
354430	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ NORBERTO				2016-06-25T BRONX		LECTURER ACTIVE	56939 per Annum	219.28	47938.44	0	0	0	0	
354431	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ ORLANDO J				1987-07-01T BRONX		PROFESSOR ACTIVE	128485 per Annum	260.72	143826.28	0	0	0	1000	
354432	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ PEDRO				2013-07-02T BRONX		CAMPUS SEC ACTIVE	33741 per Annum	2045.72	40149.64	315.75	9849.54	0	5887.37	
354433	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ ORLANDO J				1977-09-01T BRONX		PROFESSOR SEASONAL	79.15 per Day	0	52.63	0	0	0	0	
354434	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ PEGGY M				2013-09-20T BRONX		COLLEGE ASI CEASED	10.37 per Hour	0	385.9	0	0	0	0.71	
354435	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ NEIL V				2010-02-02T BRONX		DISTINGUISH CEASED	104050 per Annum	0	1772.01	0	0	0	0	
354436	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ ORLANDO J				1977-09-01T BRONX		NON-TEACH ACTIVE	53.82 per Hour	15	807.3	0	0	0	0	
354437	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ PETE P				1993-04-30T BRONX		LOCKSMITH ACTIVE	225.52 per Day	2085.72	60890.4	3.75	158.57	0	3360.37	
354438	2017	COMMUNITY COLLEGE (HOSTOS HERRELL TODD M				2007-06-04T BRONX		ADJUNCT LEI CEASED	37.92 per Day	0	106.65	0	0	0	0	
354439	2017	COMMUNITY COLLEGE (HOSTOS HEREDIA ROSA				2015-06-24T BRONX		ADJUNCT AS SEASONAL	151.66 per Day	160	27823.16	0	0	0	697.2	
354440	2017	COMMUNITY COLLEGE (HOSTOS HERNAIZ FEUX B				2015-06-12T BRONX		CAMPUS PE ACTIVE	38791 per Annum	2077.72	39172.56	387.75	12726.48	0	5448.02	
354441	2017	COMMUNITY COLLEGE (HOSTOS HEREDIA RAFAEL				2012-01-31T BRONX		LECTURE CEASED	59177 per Annum	0	1030.51	0	0	0	0	
354442	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ CLAUDIA F				2005-10-31T BRONX		CUNY ADMIR CEASED	45265 per Annum	0	3146.62	0	4.3	0	0	
354443	2017	COMMUNITY COLLEGE (HOSTOS HERAS HENRY P				2015-10-30T BRONX		CONTINUING ACTIVE	36.64 per Hour	0	50.83	0	0	0	0	
354444	2017	COMMUNITY COLLEGE (HOSTOS HEREDIA RAFAEL				2013-01-31T BRONX		ADJUNCT LEI CEASED	85.12 per Day	0	75.84	0	0	0	0	
354445	2017	COMMUNITY COLLEGE (HOSTOS HERNANDEZ CLAUDIA F				2014-06-02T BRONX		ASSISTANT TA ACTIVE	55837 per Annum	1822	62840.43	0	0	0	1000	

Dataset screenshot - 1

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	
1	Fiscal Year	Payroll Num	Agency Name	Last Name	First Name	Mid Init	Agency Start Date	Work Location Borough	Title Description	Leave Status	Base Salary	Pay Basis	Regular Hours	Regular Gross Paid	OT Hours	Total OT Paid	Total Other Pay
2	2020	67	ADMIN FOR CHILDREN'S SVCS	MAKRHINSKY, IRINA			2011-10-03T00:00:00.000	MANHATTAN	STAFF ANALYST,ACTIVE	81509.00	per Annum	1820.00	97329.44	44.00	0.00	0.00	0
3	2020	67	ADMIN FOR CHILDREN'S SVCS	ACEVEDO, MIGUELINA P			2006-02-27T00:00:00.000	BRONX	CHILD PROTECTIVE SPECIALIST,SUPERVISOR,ACTIVE	86096.00	per Annum	1820.00	89599.83	142.25	9830.52	4501.92	
4	2020	67	ADMIN FOR CHILDREN'S SVCS	ARGUETA, REGINA A			2015-02-09T00:00:00.000	BROOKLYN	CHILD PROTECTIVE SPECIALIST,ACTIVE	60327.00	per Annum	1820.00	58937.6	336	14387.01	2027.07	
5	2020	67	ADMIN FOR CHILDREN'S SVCS	SANTOS, WESLEY			2017-02-21T00:00:00.000	MANHATTAN	PROGRAM E ACTIVE	82987.00	per Annum	1820.00	79325.99	0	0	0	
6	2020	67	ADMIN FOR CHILDREN'S SVCS	BURKE, KEVIN	C		2000-03-13T00:00:00.000	MANHATTAN	ADMINISTRATIVE ACTIVE	90764.00	per Annum	1820.00	89370.34	0	0	795.64	
7	2020	67	ADMIN FOR CHILDREN'S SVCS	SHAH, HETAL M			2017-06-12T00:00:00.000	BRONX	AGENCY ATT ATTIVE	91563.00	per Annum	1820.00	90972.86	147.25	9746.24	3196.65	
8	2020	67	ADMIN FOR CHILDREN'S SVCS	DELROW, ANGELE			1997-09-22T00:00:00.000	MANHATTAN	PRINCIPAL A ACTIVE	68688.00	per Annum	1820.00	69826.26	15	927.94	1721.4	
9	2020	67	ADMIN FOR CHILDREN'S SVCS	DAY, JENNIFER			2000-02-27T00:00:00.000	QUEENS	CHILD PROTECTIVE SPECIALIST,ACTIVE	66037.00	per Annum	1820.00	60725.0	122.25	5172.59	4137.62	
10	2020	67	ADMIN FOR CHILDREN'S SVCS	WILLIAMS, INNA			1996-02-23T00:00:00.000	QUEENS	CHILD WELF ACTIVE	72525.00	per Annum	1820.00	70352.8	6.75	277.81	4784.07	
11	2020	67	ADMIN FOR CHILDREN'S SVCS	MARCEL, INGRID			1995-02-01T00:00:00.000	BRONX	CHILD PROTECTIVE SPECIALIST,ACTIVE	68037.00	per Annum	1820.00	68827.37	351.25	15379.97	4041.44	
12	2020	67	ADMIN FOR CHILDREN'S SVCS	ANGUS, DESMOND			2006-09-04T00:00:00.000	MANHATTAN	ADMINISTRATIVE ACTIVE	85031.00	per Annum	1820.00	83725.26	0	0	803.89	
13	2020	67	ADMIN FOR CHILDREN'S SVCS	HICKS, LILIAN L			1996-06-23T00:00:00.000	BRONX	CHILD PROTECITIVE SPECIALIST,ACTIVE	86419.00	per Annum	1820.00	85084.25	0	0	5004.94	
14	2020	67	ADMIN FOR CHILDREN'S SVCS	SPRINGER, ANDRE R			1996-06-23T00:00:00.000	MANHATTAN	CHILD PROTECITIVE SPECIALIST,ACTIVE	92351.00	per Annum	1820.00	99933.0	0	0	4830.4	
15	2020	67	ADMIN FOR CHILDREN'S SVCS	ABRAHAM, YAVOLA K			2014-12-02T00:00:00.000	MANHATTAN	AGENCY ATT ACTIVE	91563.00	per Annum	1820.00	94887.78	74	4412.59	4582.81	
16	2020	67	ADMIN FOR CHILDREN'S SVCS	MORALES-IS JUANA			2011-07-04T00:00:00.000	QUEENS	SUPERVISOR,ACTIVE	55853.00	per Annum	2080.00	55216.73	338.5	13875.27	2177.25	
17	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA L			2012-06-18T00:00:00.000	BRONX	YOUTH DEV/E ACTIVE	59995.00	per Annum	1750.00	58410.11	158	7351.05	6690.81	
18	2020	67	ADMIN FOR CHILDREN'S SVCS	PINNICK, MONITE N			2017-10-16T00:00:00.000	MANHATTAN	CHILD PROTECITIVE SPECIALIST,ACTIVE	60327.00	per Annum	1820.00	58929.41	289.5	12778.57	3626.93	
19	2020	67	ADMIN FOR CHILDREN'S SVCS	ANTOINE, LUCY			2006-09-04T00:00:00.000	MANHATTAN	ADMINISTRATIVE ACTIVE	80106.00	per Annum	1820.00	72346.34	0	0	2998.15	
20	2020	67	ADMIN FOR CHILDREN'S SVCS	WOODS, ANNJEANTEEN			2015-07-06T00:00:00.000	BRONX	COMMUNITY ACTIVE	65896.00	per Annum	1820.00	64876.12	0	0	1324.26	
21	2020	67	ADMIN FOR CHILDREN'S SVCS	RODOLFO, SHAVON E			2017-02-13T00:00:00.000	BRONX	COMMUNITY ACTIVE	62215.00	per Annum	1820.00	61259.83	312	13804.06	313.5	
22	2020	67	ADMIN FOR CHILDREN'S SVCS	SMITH-AGB/JUNE A			2005-08-01T00:00:00.000	MANHATTAN	CHILD PROTECITIVE SPECIALIST,ACTIVE	67392.00	per Annum	1820.00	62045.47	221	9472.76	4376.44	
23	2020	67	ADMIN FOR CHILDREN'S SVCS	HERRERA, MANUEL			1996-06-23T00:00:00.000	BRONX	CHILD PROTECITIVE SPECIALIST,ACTIVE	64727.00	per Annum	1820.00	63261.96	258.25	11029.57	4968.41	
24	2020	67	ADMIN FOR CHILDREN'S SVCS	OLIVER, DALE			1998-09-08T00:00:00.000	MANHATTAN	ASSOCIATE S ACTIVE	77083.00	per Annum	1820.00	76666.4	0	0	5437.45	
25	2020	67	ADMIN FOR CHILDREN'S SVCS	VIDAL, SIGRY E			2015-02-09T00:00:00.000	BRONX	CHILD PROTECITIVE SPECIALIST,ACTIVE	60327.00	per Annum	1820.00	58937.6	543.75	23504.42	1218.84	
26	2020	67	ADMIN FOR CHILDREN'S SVCS	KINARD, TRAMAINIE T			2016-07-18T00:00:00.000	MANHATTAN	MOTOR VEH ACTIVE	49927.00	per Annum	2080.00	49103.96	1054.5	37436.55	1239.55	
27	2020	67	ADMIN FOR CHILDREN'S SVCS	WALKER, SERENA			1996-06-23T00:00:00.000	MANHATTAN	CHILD WELF ACTIVE	86452.00	per Annum	1820.00	85116.83	3	150	4844.15	
28	2020	67	ADMIN FOR CHILDREN'S SVCS	SUAREZ, JOSEPHINE			2001-05-21T00:00:00.000	MANHATTAN	DIRECTOR O ACTIVE	95831.00	per Annum	1820.00	94375.02	0	0	0	
29	2020	67	ADMIN FOR CHILDREN'S SVCS	HAMMONDS, MONIQUE A			2017-02-13T00:00:00.000	MANHATTAN	ADMINISTRATIVE ACTIVE	94191.00	per Annum	1820.00	92213.57	0	0	3455.34	
30	2020	67	ADMIN FOR CHILDREN'S SVCS	LACEWELL, MARTIN			1997-01-06T00:00:00.000	MANHATTAN	CONGRESS G ACTIVE	63691.00	per Annum	1820.00	62705.3	540.25	25184.16	7884.17	
31	2020	67	ADMIN FOR CHILDREN'S SVCS	DOUGLAS-AI SHERRILL A			2004-01-26T00:00:00.000	BRONX	CHILD PROTECITIVE SPECIALIST,ACTIVE	86096.00	per Annum	1820.00	84774.03	244	14085.24	5014.25	
32	2020	67	ADMIN FOR CHILDREN'S SVCS	WILLIAMS, DAMITA D			1996-06-23T00:00:00.000	MANHATTAN	PRINCIPAL A ACTIVE	65758.00	per Annum	1820.00	67042.68	0	0	4219.81	
33	2020	67	ADMIN FOR CHILDREN'S SVCS	WILLIAMS-M HELENA C			2007-07-16T00:00:00.000	MANHATTAN	CHILD PROTECITIVE SPECIALIST,ACTIVE	86096.00	per Annum	1820.00	88842.54	136.75	8215.65	4227.38	
34	2020	67	ADMIN FOR CHILDREN'S SVCS	SUBAIRA, BASIRAT F			2015-04-20T00:00:00.000	MANHATTAN	STRATEGIC I ACTIVE	92700.00	per Annum	1820.00	91291.4	39.25	2173.67	117.13	
35	2020	67	ADMIN FOR CHILDREN'S SVCS	SMITH, MAI T			2018-06-25T00:00:00.000	MANHATTAN	CONGRES G ACTIVE	42731.00	per Annum	1546.00	36027.04	1021.73	34038.5	2369.59	
36	2020	67	ADMIN FOR CHILDREN'S SVCS	EVANS, CATHERINE			2008-08-01T00:00:00.000	MANHATTAN	CHILD PROTECTIVE SPECIALIST,ACTIVE	92351.00	per Annum	1820.00	88888.79	100.00	1212.00	4830.40	
37	2020	67	ADMIN FOR CHILDREN'S SVCS	WILLIAMS, ARTHUR			2008-03-27T00:00:00.000	MANHATTAN	CHILD PROTECTIVE SPECIALIST,ACTIVE	86096.00	per Annum	1820.00	85521.0	100.00	0.00	0.00	
38	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
39	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
40	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
41	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
42	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
43	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
44	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
45	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
46	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
47	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
48	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
49	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
50	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
51	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
52	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
53	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
54	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
55	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
56	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
57	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81	
58	2020	67	ADMIN FOR CHILDREN'S SVCS	WHITE, QUANICUA			2012-06-18T00:00:00.000	BROOKLYN	YOUTH DEVELOPMENT SPECIALIST,ACTIVE	59995.00	per Annum	1750.00	58410.74	151.00	7351.8	6590.81</td	

## Raw data screenshot - 2

## Dataset after MapReduce

BRONX, 2015, POLICE DEPARTMENT, 357.77, 22147.74, 108875.49  
BRONX, 2018, POLICE DEPARTMENT, 59.13, 2109.36, 55327.06  
BRONX, 2016, POLICE DEPARTMENT, 41.25, 2817.61, 61492.44  
BRONX, 2017, POLICE DEPARTMENT, 242.25, 22182.34, 124892.769999999999  
BRONX, 2017, POLICE DEPARTMENT, 88.17, 3499.45, 56625.009999999995  
BRONX, 2017, POLICE DEPARTMENT, 195.08, 18961.07, 145875.09  
BRONX, 2016, POLICE DEPARTMENT, 384.83, 21720.32, 114479.08  
BRONX, 2020, POLICE DEPARTMENT, 240.67, 8658.92, 111157.17  
BRONX, 2019, POLICE DEPARTMENT, 30.30, 2278.18, 108282.819999999999  
BRONX, 2015, POLICE DEPARTMENT, 0, 0.00, 24342.350000000002  
BRONX, 2016, POLICE DEPARTMENT, 549, 50522.10, 174474.04  
BRONX, 2017, POLICE DEPARTMENT, 142.75, 5432.70, 55840.32  
BRONX, 2018, POLICE DEPARTMENT, 41, 2123.08, 64925.6  
BRONX, 2015, POLICE DEPARTMENT, 308, 9279.39, 49619.11  
BRONX, 2017, POLICE DEPARTMENT, 20, 805.93, 56944.47  
BRONX, 2019, POLICE DEPARTMENT, 42.33, 3024.28, 103764.65  
BRONX, 2020, POLICE DEPARTMENT, 122.00, 3109.95, 53304.659999999996  
BRONX, 2017, POLICE DEPARTMENT, 224.48, 11478.42, 70140.38  
BRONX, 2018, POLICE DEPARTMENT, 29.5, 2126.31, 100009.959999999999  
BRONX, 2017, POLICE DEPARTMENT, 95.45, 5113.75, 67561.41  
BRONX, 2016, POLICE DEPARTMENT, 410, 25508.29, 126601.2300000001  
BRONX, 2017, POLICE DEPARTMENT, 30, 3901.00, 110196.22  
BRONX, 2017, POLICE DEPARTMENT, 201.83, 25004.83, 131555.32  
BRONX, 2016, POLICE DEPARTMENT, 532.78, 39022.50, 136057.22  
BRONX, 2015, POLICE DEPARTMENT, 400.18, 23749.15, 118730.98  
BRONX, 2019, POLICE DEPARTMENT, 33.42, 1346.43, 68592.81  
BRONX, 2015, POLICE DEPARTMENT, 347.67, 20455.85, 113260.70000000001  
BRONX, 2015, POLICE DEPARTMENT, 565.25, 21371.87, 78423.33  
BRONX, 2020, POLICE DEPARTMENT, 520.22, 29395.57, 139531.83000000002  
BRONX, 2019, POLICE DEPARTMENT, 0.00, 0.00, 188989.5  
BRONX, 2018, POLICE DEPARTMENT, 28.25, 2140.12, 110443.66  
BRONX, 2017, POLICE DEPARTMENT, 293.75, 7298.56, 48346.42  
BRONX, 2020, POLICE DEPARTMENT, 355.25, 10910.68, 54273.32  
BRONX, 2019, POLICE DEPARTMENT, 237.42, 13751.01, 115215.64  
BRONX, 2019, POLICE DEPARTMENT, 0.00, 0.00, 25938.0  
BRONX, 2019, POLICE DEPARTMENT, 246.10, 9949.98, 67661.76  
BRONX, 2020, POLICE DEPARTMENT, 614.92, 26872.99, 125298.71  
BRONX, 2016, POLICE DEPARTMENT, 49.25, 1265.73, 42275.77000000004  
BRONX, 2018, POLICE DEPARTMENT, 47.67, 2145.93, 70392.64  
BRONX, 2016, POLICE DEPARTMENT, 72.68, 1861.30, 47035.82  
BRONX, 2018, POLICE DEPARTMENT, 61.22, 1983.19, 48660.5  
BRONX, 2019, POLICE DEPARTMENT, 213.75, 15592.27, 119816.0400000001  
BRONX, 2015, POLICE DEPARTMENT, 345.68, 21272.80, 121778.52  
BRONX, 2016, POLICE DEPARTMENT, 511.42, 36732.95, 140203.75  
BRONX, 2019, POLICE DEPARTMENT, 464.25, 15867.81, 63988.49  
BRONX, 2019, POLICE DEPARTMENT, 50.58, 5271.81, 186014.44  
BRONX, 2017, POLICE DEPARTMENT, 0, 0.00, 44178.89  
BRONX, 2019, POLICE DEPARTMENT, 8.42, 279.57, 44408.159999999996  
BRONX, 2017, POLICE DEPARTMENT, 113.75, 4556.84, 66972.93  
BRONX, 2020, POLICE DEPARTMENT, 268.25, 15358.02, 118756.59  
BRONX, 2019, POLICE DEPARTMENT, 233.82, 17715.25, 119467.959999999999  
BRONX, 2015, POLICE DEPARTMENT, 146.92, 9582.60, 100893.90000000001  
BRONX, 2020, POLICE DEPARTMENT, 140.40, 8781.16, 108618.27  
BRONX, 2018, POLICE DEPARTMENT, 27.57, 2156.79, 110671.799999999999  
BRONX, 2016, POLICE DEPARTMENT, 379.38, 27216.78, 130983.34  
BRONX, 2017, POLICE DEPARTMENT, 298.87, 25334.74, 129738.35  
BRONX, 2018, POLICE DEPARTMENT, 60.17, 2158.30, 55708.36  
BRONX, 2016, POLICE DEPARTMENT, 387.42, 30279.88, 151772.68  
BRONX, 2015, POLICE DEPARTMENT, 647.15, 43505.18, 142977.75  
BRONX, 2020, POLICE DEPARTMENT, 411.42, 22875.56, 139479.55  
BRONX, 2019, POLICE DEPARTMENT, 66.75, 5664.00, 107116.33  
BRONX, 2015, POLICE DEPARTMENT, 214.5, 6114.69, 56636.05  
BRONX, 2017, POLICE DEPARTMENT, 58.8, 4750.65, 106889.15  
BRONX, 2020, POLICE DEPARTMENT, 352.58, 20591.34, 125301.269999999999  
BRONX, 2018, POLICE DEPARTMENT, 46.17, 2168.74, 74345.05  
BRONX, 2017, POLICE DEPARTMENT, 478.97, 42657.25, 163718.01  
BRONX, 2020, POLICE DEPARTMENT, 249.50, 5775.95, 52991.45  
BRONX, 2017, POLICE DEPARTMENT, 561.17, 49712.68, 159316.88  
BRONX, 2020, POLICE DEPARTMENT, 254.50, 9436.53, 61924.06

### Preprocessed data screenshot - 1

BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,71,2802.47,72512.23  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,1.25,40.41,58596.020000000004  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,112095.39  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,8,499.38,96697.26000000001  
 BROOKLYN,2017,DEPT OF ENVIRONMENT PROTECTION,404.25,26948.00,96879.46  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,198.5,7306.53,71015.27  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,128.00,3864.66,58988.53999999999  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,122,2295.39,33165.12  
 BROOKLYN,2019,DEPT OF HEALTH/MENTAL HYGIENE,0.00,0.00,66064.56  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,536.75,15838.90,58988.740000000005  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,4.88,66653.68  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,19,406.43,40825.68  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,67,2470.36,50178.84  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,43176.41  
 BROOKLYN,2019,DEPT OF HEALTH/MENTAL HYGIENE,0.00,0.00,66681.88  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,59284.21  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,9,173.11,33152.560000000005  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,29.75,1776.75,98949.16  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,29.5,875.32,43432.58  
 BROOKLYN,2020,DEPT. OF HOMELESS SERVICES,5.75,191.09,62415.17999999999  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,0.00,0.00,35483.26  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,42499.92  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,113.25,5739.79,68685.54000000001  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,8.5,276.17,62429.659999999996  
 BROOKLYN,2016,NYC HOUSING AUTHORITY,5,81.91,31106.91  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,50.61,42586.03  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,63345.0  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,175.5,5261.59,49973.82999999994  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,243.25,6099.55,43151.880000000005  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,317.25,10417.00,55938.88  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,116.75,3004.65,35979.5  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,279.25,8997.43,55421.4  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,38.75,1446.17,53360.82  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,7,243.69,53061.92  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,23.5,850.65,57483.11  
 BROOKLYN,2019,DEPT OF HEALTH/MENTAL HYGIENE,0.00,0.00,66151.56  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,359.5,11071.48,57734.03  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,222.75,5633.09,50256.14  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,48995.17  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,5,187.94,69148.04000000001  
 BROOKLYN,2017,DEPT OF ENVIRONMENT PROTECTION,512.5,34549.56,129859.55999999998  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,124.00,3429.18,40452.53  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,342,24105.00,135979.38  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,165.25,3925.06,41114.52  
 BROOKLYN,2019,DEPT OF HEALTH/MENTAL HYGIENE,0.00,1.88,66153.54  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,49070.17  
 BROOKLYN,2019,DEPT OF HEALTH/MENTAL HYGIENE,0.00,0.00,66183.79  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,150.00,3705.46,40931.54999999996  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,248,11483.75,71835.85  
 BROOKLYN,2016,NYC HOUSING AUTHORITY,183.75,7857.76,64237.78000000006  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,22.75,1061.86,47319.44  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,120.00,3199.77,40558.38  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,141,6973.24,83282.12000000001  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,858.75,34472.66,87898.72  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,133.50,3679.37,40728.54  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,181.50,4547.58,36622.72  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,0.00,0.00,41408.58  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,40411.0  
 BROOKLYN,2020,DEPT. OF HOMELESS SERVICES,716.25,19235.10,61772.65  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,9.25,275.81,53127.78999999999  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,1,25.18,58309.83  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,7.93,56858.78  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,58.5,2039.04,54807.77000000004  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,172.00,4154.48,35839.82999999994  
 BROOKLYN,2020,DEPT. OF HOMELESS SERVICES,0.00,0.00,62215.0  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,0,0.00,43639.77000000004  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,125.75,4717.61,61428.23  
 BROOKLYN,2017,HRA/DEPT OF SOCIAL SERVICES,399.5,15373.05,74903.91  
 BROOKLYN,2019,NYC HOUSING AUTHORITY,265.25,6833.54,44327.32

Preprocessed data screenshot - 2

Preprocessed data screenshot - 3

QUEENS,2018,POLICE DEPARTMENT,399.02,30095.05,129250.0  
 QUEENS,2018,POLICE DEPARTMENT,22.38,591.58,45173.64  
 QUEENS,2020,POLICE DEPARTMENT,967.20,36385.40,143193.9499999998  
 QUEENS,2020,POLICE DEPARTMENT,194.37,9070.63,143319.17  
 QUEENS,2018,POLICE DEPARTMENT,412.83,30103.90,130333.51999999999  
 QUEENS,2020,POLICE DEPARTMENT,697.55,36999.23,143177.83000000002  
 QUEENS,2020,POLICE DEPARTMENT,621.90,31683.65,143241.49  
 QUEENS,2015,POLICE DEPARTMENT,0.0.00,38997.95  
 QUEENS,2018,POLICE DEPARTMENT,282,30120.76,132074.03999999998  
 QUEENS,2016,POLICE DEPARTMENT,144.67,7571.65,73690.38  
 QUEENS,2018,POLICE DEPARTMENT,24,591.95,36090.89  
 QUEENS,2020,POLICE DEPARTMENT,308.45,15859.98,143472.91  
 QUEENS,2016,POLICE DEPARTMENT,540,25228.14,97483.97  
 QUEENS,2016,POLICE DEPARTMENT,92.48,5368.19,89755.44  
 QUEENS,2016,POLICE DEPARTMENT,314.08,22611.10,115859.21  
 QUEENS,2016,POLICE DEPARTMENT,224.92,9585.21,68779.03  
 QUEENS,2018,POLICE DEPARTMENT,440.5,30134.86,128028.3  
 QUEENS,2020,POLICE DEPARTMENT,857.08,37829.19,143160.75  
 QUEENS,2018,POLICE DEPARTMENT,21.5,596.58,39430.63000000005  
 QUEENS,2016,POLICE DEPARTMENT,0.10.25,47890.43  
 QUEENS,2018,POLICE DEPARTMENT,10.5,598.66,39327.26  
 QUEENS,2018,POLICE DEPARTMENT,17.75,613.87,51862.53  
 QUEENS,2016,POLICE DEPARTMENT,316.83,20853.96,110273.13999999998  
 QUEENS,2018,POLICE DEPARTMENT,11.58,614.63,68430.39  
 QUEENS,2018,POLICE DEPARTMENT,19.77,617.08,43893.5  
 QUEENS,2015,POLICE DEPARTMENT,241,16529.61,121029.25  
 QUEENS,2018,POLICE DEPARTMENT,905,30155.59,77722.79  
 QUEENS,2018,POLICE DEPARTMENT,315.92,30163.35,173347.99  
 QUEENS,2018,POLICE DEPARTMENT,19.25,618.11,45178.85  
 QUEENS,2018,POLICE DEPARTMENT,0.624.81,55349.28999999999  
 QUEENS,2016,POLICE DEPARTMENT,131.85,10535.65,131276.35  
 QUEENS,2018,POLICE DEPARTMENT,408.93,30180.00,129712.47  
 QUEENS,2015,TAXI & LIMOUSINE COMMISSION,121,3702.21,46759.21  
 QUEENS,2018,POLICE DEPARTMENT,22,626.18,40381.08  
 QUEENS,2020,POLICE DEPARTMENT,500.08,27314.39,143217.52  
 QUEENS,2020,POLICE DEPARTMENT,488.12,27662.61,143215.49  
 QUEENS,2015,POLICE DEPARTMENT,1,17.56,34984.74  
 QUEENS,2018,POLICE DEPARTMENT,19.17,628.65,45200.05  
 QUEENS,2018,POLICE DEPARTMENT,382.9,30195.65,138956.53999999998  
 QUEENS,2018,POLICE DEPARTMENT,16,629.19,48581.36  
 QUEENS,2018,POLICE DEPARTMENT,401.82,30201.46,131372.72  
 QUEENS,2018,POLICE DEPARTMENT,8.58,630.35,101619.20000000001  
 QUEENS,2015,TAXI & LIMOUSINE COMMISSION,0.0.56,30524.56  
 QUEENS,2015,POLICE DEPARTMENT,628,16519.16,54528.02000000004  
 QUEENS,2015,TAXI & LIMOUSINE COMMISSION,14.25,555.45,66034.55  
 QUEENS,2015,POLICE DEPARTMENT,278.45,22405.39,146730.34  
 QUEENS,2018,POLICE DEPARTMENT,937,30205.69,75315.35  
 QUEENS,2018,POLICE DEPARTMENT,11,637.67,81581.25  
 QUEENS,2020,POLICE DEPARTMENT,600.50,27320.34,143185.81  
 QUEENS,2018,POLICE DEPARTMENT,8.58,637.85,102540.36  
 QUEENS,2020,POLICE DEPARTMENT,105.33,2560.09,62841.03999999999  
 QUEENS,2018,POLICE DEPARTMENT,439.25,30217.09,128375.29  
 QUEENS,2018,POLICE DEPARTMENT,8.5,639.08,106169.17  
 QUEENS,2018,POLICE DEPARTMENT,299.83,30223.07,171864.03  
 QUEENS,2018,POLICE DEPARTMENT,393.82,30225.42,129972.44  
 QUEENS,2015,POLICE DEPARTMENT,180.07,17982.45,139173.6  
 QUEENS,2020,POLICE DEPARTMENT,172.50,6332.19,57612.75  
 QUEENS,2016,POLICE DEPARTMENT,406.57,34556.44,154143.76  
 QUEENS,2020,POLICE DEPARTMENT,74.75,4903.83,154807.5  
 QUEENS,2018,POLICE DEPARTMENT,23.5,642.82,39352.92  
 QUEENS,2018,POLICE DEPARTMENT,20.33,650.76,45579.48  
 QUEENS,2020,TAXI & LIMOUSINE COMMISSION,156.75,5238.30,61847.92000000006  
 QUEENS,2018,POLICE DEPARTMENT,962,30251.11,74821.09  
 QUEENS,2018,POLICE DEPARTMENT,17.92,655.42,56888.79  
 QUEENS,2016,POLICE DEPARTMENT,109.5,8139.88,112176.83  
 QUEENS,2020,POLICE DEPARTMENT,87.50,3517.86,59219.91  
 QUEENS,2018,POLICE DEPARTMENT,281.75,30263.72,181992.48  
 QUEENS,2015,POLICE DEPARTMENT,578,34511.40,102425.4  
 QUEENS,2018,POLICE DEPARTMENT,397.3,30264.67,129696.28  
 QUEENS,2016,POLICE DEPARTMENT,480.9,35035.84,141078.47999999998

Preprocessed data screenshot - 4

# Shell commands

---

## Running commands:

```
# Make the project directory
hadoop fs -mkdir project

# Setting hadoop classpath and creating the jar
javac -classpath `hadoop classpath` PayrollFiscalMapper.java
javac -classpath `hadoop classpath` PayrollFiscalReducer.java
javac -classpath `hadoop classpath`:. PayrollFiscal.java
jar cvf PayrollFiscal.jar *.class

# Putting the data file onto hadoop
hadoop fs -put citywide-payroll-data-fiscal-year.csv project/

# Command to run the task
hadoop jar PayrollFiscal.jar PayrollFiscal project/citywide-payroll-data-fiscal-year.c
sv project/output7

# Getting the output to dataproc
hadoop fs -get project/output7
```

## Other shell commands:

```
# Log in to Google cloud
gcloud compute ssh nyu-dataproc-m --project hpc-dataproc-19b8 --zone us-central1-f

# scp command to push data from local
gcloud compute scp --recurse /Users/rachitjain/Desktop/Courses/Sem3/RBDA/Assignments/P
roject/nyc-crime-arrest-payroll rj2219_nyu_edu@nyu-dataproc-m:project/ --project hpc-d
ataproc-19b8 --zone us-central1-f

# scp command to get the output to local
gcloud compute scp --recurse rj2219_nyu_edu@nyu-dataproc-m:project/nyc-crime-arrest-pa
yroll/output7/ /Users/rachitjain/Desktop/Courses/Sem3/RBDA/Assignments/Project/nyc-cri
me-arrest-payroll/results/ --project hpc-dataproc-19b8 --zone us-central1-f

# Extra commands
hadoop fs -ls project/
hadoop fs -ls project/output7
hadoop fs -cat project/output7/part-r-00000

# Count how many records left after preprocessing
wc -l part-r-00000
```

## Team Members

---

- [1] Jahnvi Arya (ja4158)
- [2] Rachit Jain (rj2219)
- [3] Shobhit Sinha (ss13881)

## References

---

- [1] York, C. of N. (2021, January 1). *NY citywide payroll data (fiscal year)*. Kaggle. Retrieved November 27, 2022, from <https://www.kaggle.com/datasets/new-york-city/ny-citywide-payroll-data-fiscal-year>
- [2] <https://github.com/rachitjain2706/nyc-crime-arrest-payroll>
- [3] <https://www.indeed.com/career-advice/career-development/work-weeks-in-year>
- [4] <https://www.ny.gov/new-york-states-minimum-wage/new-york-states-minimum-wage#:~:text=The%20Minimum%20Wage%20Act>