

# PROGRESS REPORT 2 – 12/14/2017

# WHY ARE SORTING ALGORITHMS IMPORTANT?

- 25-50% OF WHAT COMPUTERS DO REQUIRE SORTING ALGORITHMS FOR
  - EFFICIENCY
  - SPEED
  - PROCESSING DATA
    - SEARCH ENGINES (SORTING THROUGH ALL WEBSITES TO FIND RELEVANT ONES)
    - GRAPHICAL PROCESSES (LAYERING IS NECESSARY)
    - GOVERNMENT ORGANIZATIONS, FINANCIAL INSTITUTIONS, COMMERCIAL ENTERPRISES (TRANSACTIONS, ACCOUNTS, PROCESSING DATA)
    - EVENT DRIVEN SIMULATIONS

# WHAT IS BIG DATA?

- BIG DATA
  - COLLECTED BY GOVERNMENT, NASA, RESEARCH INSTITUTIONS
  - DATA SETS SO LARGE AND COMPLEX THAT IT SYSTEMS CAN NOT HANDLE THEM- DIFFICULT TO PROCESS DATA
  - VOLUME, VELOCITY, VARIETY, AND VERACITY OF DATA EXCEED STORAGE OR COMPUTING CAPACITY
  - CHALLENGE OF INDEXING, SEARCHING, TRANSFERRING, AND SO ON ALL INCREASE EXPONENTIALLY

# SORTING BIG DATA

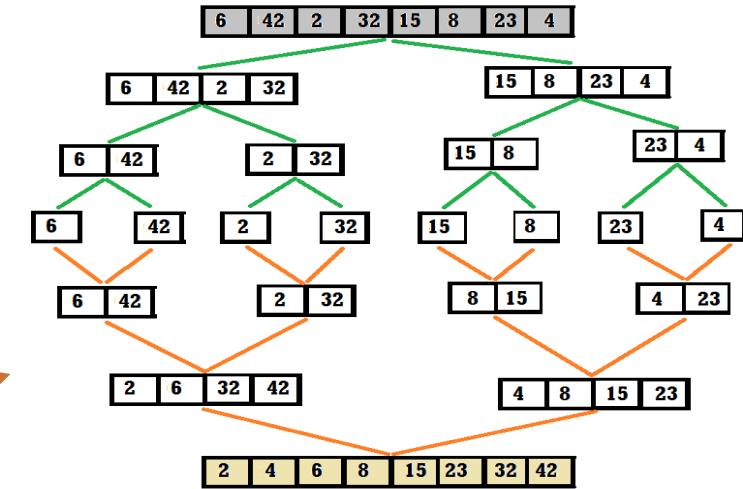
- IMAGINE YOU HAVE **123,456,789** DATA ELEMENTS  
YOU NEED TO SORT

- YOU WILL USE SORTING ALGORITHMS  
(HOPEFULLY)
- MERGE SORT = 998952457 OPERATIONS

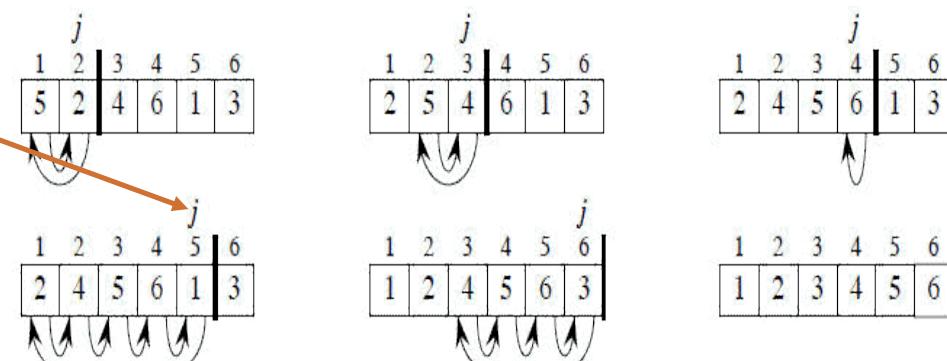
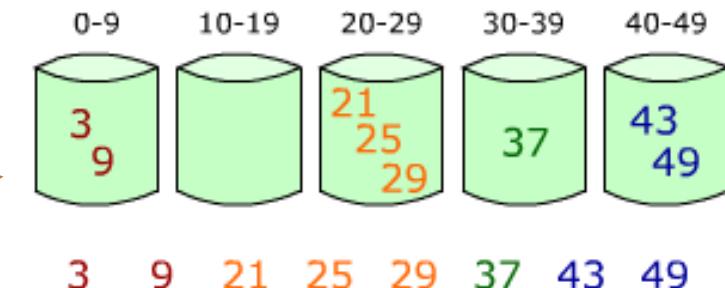
- BUCKET SORT =  $1.52 \times 10^{16}$  OPERATIONS

- INSERTION SORT =  $1.52 \times 10^{16}$  OPERATIONS

- COMPANIES HAVE MANY SERVERS AND SUPERCOMPUTERS SO THEY CAN PERFORM THIS IN MINUTES USING THE SAME ALGORITHMS



[www.Instanceof.com](http://www.Instanceof.com)



# WHY IS IT DIFFICULT?

- TOO MANY PASSES (RECURSIVE ITERATIONS)
- EACH PASS ALLOCATES ITS OWN HEAP IN THE MEMORY (EACH PASS HAS ITS OWN SPACE)
- MERGE SORT – MANY RECURSIONS
- BUCKET SORT – MEMORY IS WASTED, DISTRIBUTED DATA = UNEVEN DISTRIBUTION AMONGST BUCKETS (ALTHOUGH THIS IS A SMART IDEA)
- INSERTION SORT – SOLELY BASED ON COMPARING EACH ELEMENT TO EVERY OTHER ONE

**TOO MUCH RAM REQUIRED, TOO MANY PASSES, TOO MUCH TIME**

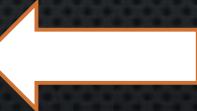
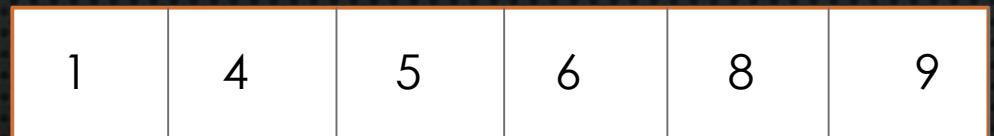
# A BETTER WAY TO SORT

- REDUCE COMPUTE TIME AND RAM REQUIREMENTS
- USE STATISTICAL ANALYSIS TO DETERMINE BUCKET RANGES BASED ON DATA SET
- <1000 ELEMENTS IN EACH BUCKET (NUMBER OF ELEMENTS ARE EQUAL)
- SHOULD BE ABLE TO SORT DATA SETS THAT FOLLOW ANY TYPE OF DISTRIBUTION
- WORST CASE OF  $O(N)$

# A VISUAL REPRESENTATION



$\mu = 5.5$   
Median = 5.5  
 $\delta = 2.88$



Bucket 1

-28 to -.58

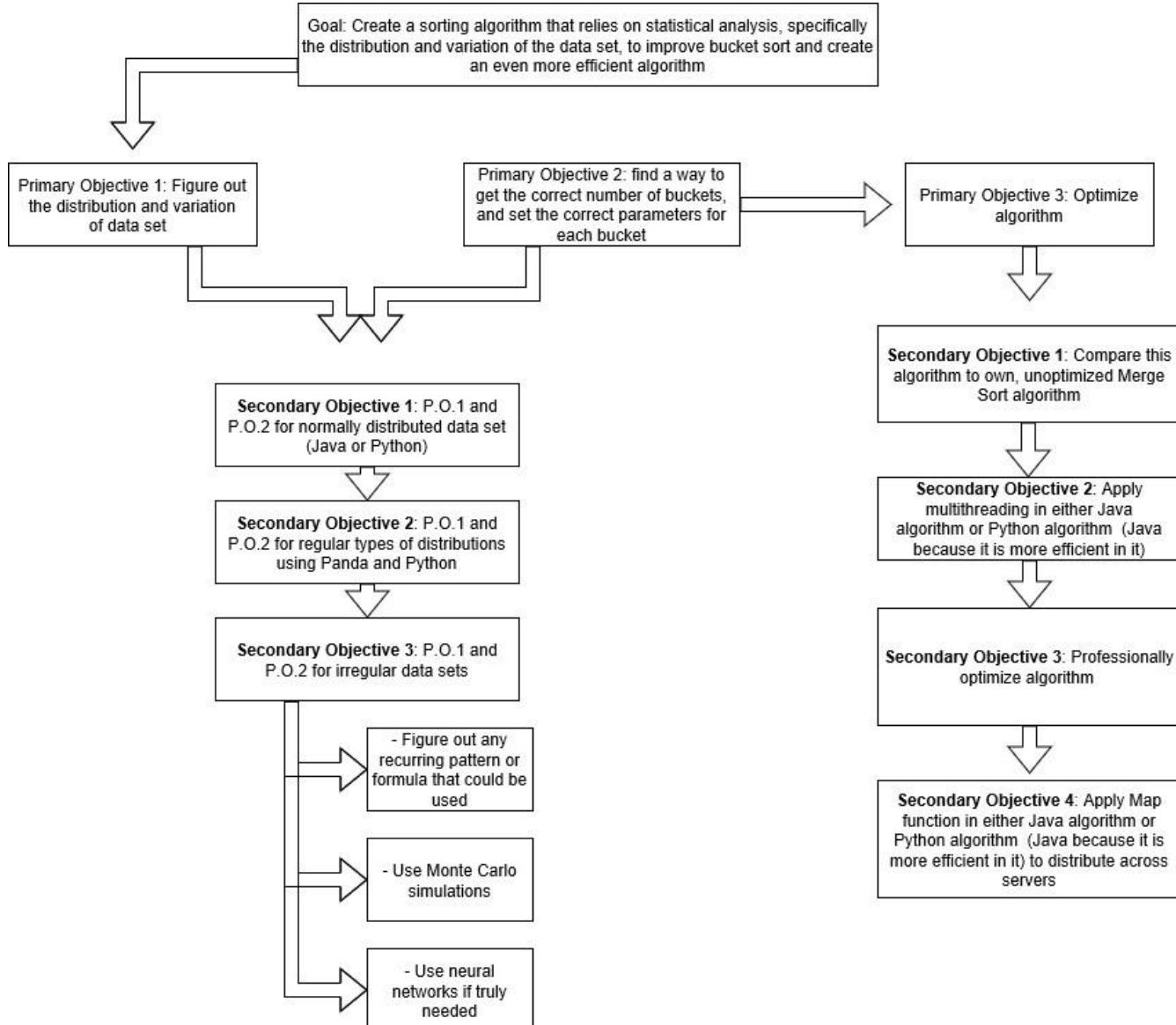
Bucket 2

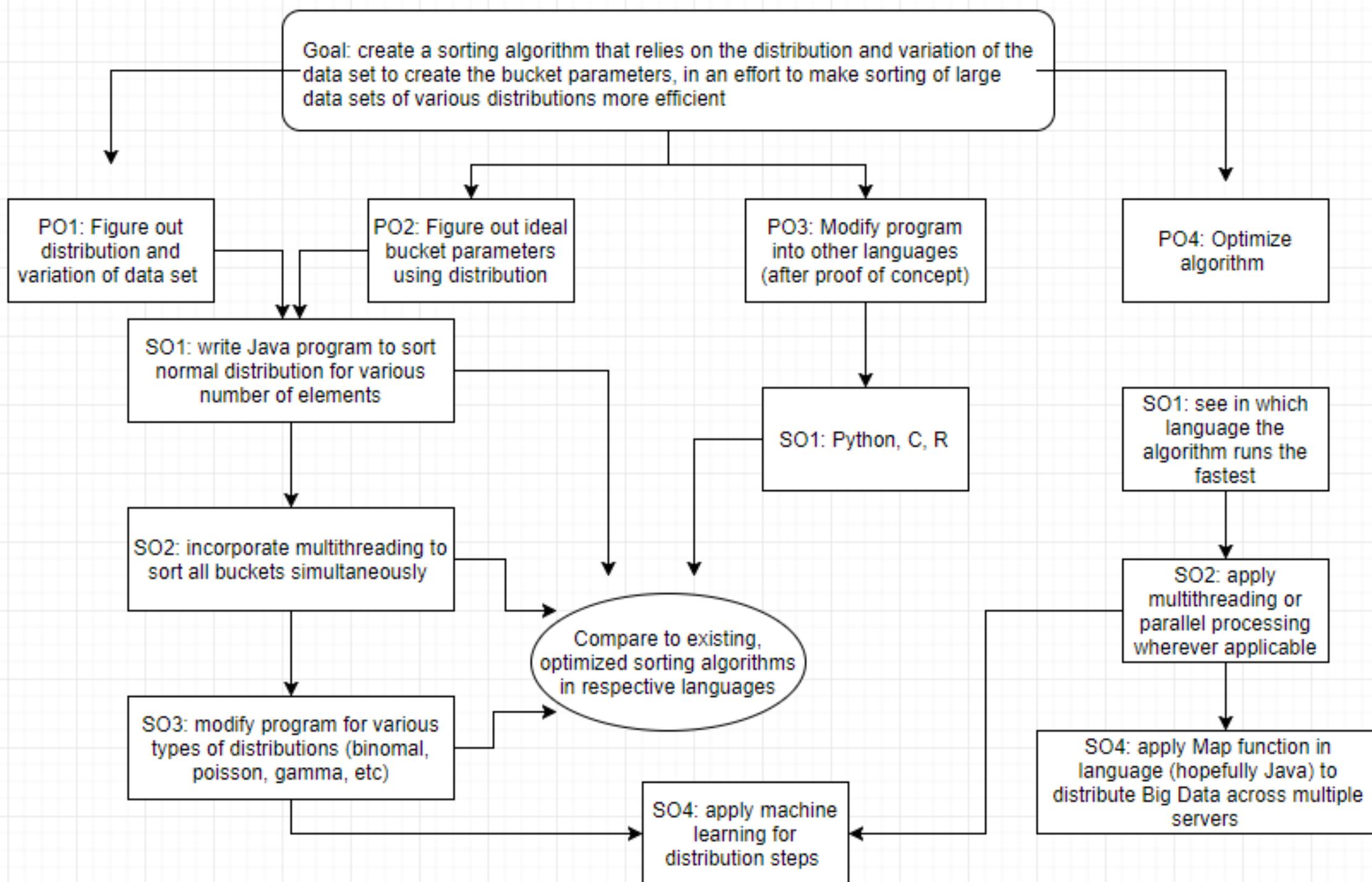
-.188 to .08

Bucket 3

0.18 to 18







# PROGRESS UPTO PROGRESS REPORT 1

- CREATED BACKBONE OF ALGORITHM IN JAVA
  - TECHNICAL ERRORS IN PRELIMINARY CODE TOOK TIME TO FIGURE OUT
- COLLECTED DATA WHEN COMPARED TO PARALLEL SORT/MERGE SORT
- TRANSFERRING CODE INTO PYTHON AND LEARNING PANDA
- STARTED COMMUNICATING WITH Ms. KUMAR

# PROGRESS UPTO NOW

- FINISHED PROGRAMMING NORMAL DISTRIBUTION CODE IN PYTHON
- ENCOUNTERED MANY PROBLEMS WHILE TRYING TO REDUCE RUN TIME
  - EMAILED Ms. KUMAR ABOUT THESE ISSUES
  - TRIED TO FIGURE OUT SOLUTIONS USING NUMPY ARRAYS AND VARIOUS NUMPY METHODS
- CHANGED PLAN DUE TO NET PROGRESS AND TIME
- LEARNING HOW TO MULTITHREAD IN JAVA
  - COMPARED CURRENT ALGORITHM TO MERGE SORT, BUCKET SORT

# PYTHON PROGRAM

```
1 import math
2 import random
3 import datetime
4 class Statsort:
5     #program to compare python embedded sort and new sort being created
6     #main - creating integer normal distribution list of varying lengths, print out times
7     def main():
8         upperlimit=10000000
9         intBucketf=[]
10        for i in range (upperlimit):
11            intBucketf.append((int)(random.uniform(0,upperlimit)))
12        Statsort.sorterSimple(intBucketf)
13        Statsort.sorterBucket(intBucketf,upperlimit)
14    def sorterSimple(intBucket):
15        print (datetime.datetime.now())
16        intBucket.sort() ← list - already optimized in python
17        print (datetime.datetime.now())
18        #sorting the integer list intBucket using the embedded sort in python
```

- Creating random normal distribution data set
- Sorting using built-in python sort (quicksort)

# PYTHON PROGRAM

```
def sorterBucket(intBucket,upperlimit):
    print (datetime.datetime.now())
    buckets=(int)(round(upperlimit/1000)) #creating the number of buckets by dividing list length (elements) by 1000 - elements in each bucket
    intbucket1=[[1000]]*buckets #creating a 2D list in which number of rows is buckets and columns is 1000
    print (datetime.datetime.now())
    Min=intBucket[0] #setting both min and max for first element in list
    Max=intBucket[0]
    sum=0
    for i in range (upperlimit):
        Min=min(Min,intBucket[i]) #comparing each value to get min and max per loop
        Max=max(Max,intBucket[i])
        sum = sum+intBucket[i] #summing all elements
    Mean=round(sum/upperlimit)
    Range=Max-Min
    print (datetime.datetime.now())
    lowerRange=Mean-Min #half of the range of the data set - used to create buckets for each range
    upperRange=Max-Mean
```

- Creating 2D list in which elements are going to be stored

- Obtaining Min, Max, Sum through use of FOR LOOP #1

# PYTHON PROGRAM

```
lowerRange=Mean-Min #half of the range of the data set - used to create buckets for each range
upperRange=Max-Mean
lowerGroups=round(lowerRange/Range) #number of buckets for lowerRange
lowerGroups=max(1,lowerGroups) #if the number of elements 1000 or less than 1000, number of groups will be 1 instead of 0
upperGroups=buckets-lowerGroups #number of buckets for upperRange
print (datetime.datetime.now())
for x in range(upperlimit):
    value=intBucket[x] #calling value at index x of list
    if Mean>value: #determining which range it will fall into
        bucketSelect=(int)(math.ceil(lowerGroups*(value-Min)/lowerRange)) #seeing where in the lowerRange does this element lie in - distance from mean
    else:
        bucketSelect=lowerGroups+(int)(math.ceil(upperGroups*(value-Mean)/upperRange)) #seeing where in the upperRange does element lie - distance from mean
    bucketSelect=(int)((value-Min)/Range)*upperlimit/1000 #reinstantiating bucketSelect
print (datetime.datetime.now())
for y in range(buckets):
    intbucket1.sort() #uses python sorting algorithm to sort buckets
print (datetime.datetime.now())
```

- Creating buckets

- Assigning elements to their respective buckets based on distance from mean

# DATA COLLECTED (IT DOES NOT WORK)

Elements	Python Sort			New Sort		
	Start Time (s)	End Time (s)	Total Time (s)	Start Time (s)	End Time (s)	Total Time (s)
1000	16.29	16.32	0.03	27.11	27.18	0.06
10000	13.76	13.82	0.06	24.13	24.25	0.12
100000	38.08	38.26	0.18	10.49	11.30	0.81
1000000	4.98	8.52	3.54	42.36	48.70	6.34
10000000	8.01	48.71	40.69	47.59	104.15	56.57

# DATA CONTINUED – WHY IS IT SLOWER THAN JAVA?

Elements	Parallel Sort (Optimized)				New Sort				Percent Faster (%)*
	Trial 1 (s)	Trial 2 (s)	Trial 3 (s)	Average (s)	Trial 1 (s)	Trial 2 (s)	Trial 3 (s)	Average (s)	
1000	0.31	0.20	0.27	0.26	0.01	0.02	0.02	0.01	95.40
10000	0.31	0.36	0.37	0.35	0.02	0.01	0.02	0.01	96.04
100000	0.24	0.25	0.28	0.26	0.05	0.08	0.07	0.07	74.03
1000000	0.62	0.68	0.51	0.61	0.16	0.16	0.13	0.15	74.82
10000000	0.48	1.25	1.24	0.99	0.26	0.54	0.47	0.42	57.32
100000000	5.71	5.91	5.82	5.81	3.31	3.06	3.40	3.26	44.01
10000000000	HEAP ERROR								

- HERE, IN JAVA, PROGRAM RUNS SIGNIFICANTLY FASTER THAN PARALLEL SORT, MERGE SORT, AND BUCKET SORT
- WHY IS THE PYTHON SIGNIFICANTLY SLOWER AND INEFFICIENT

# FURTHER INVESTIGATION

Trials	Time for Stat Collection (s)	Time for Bucketing (s)	Total Time (s)
1000	0.018	0.000	0.018
10000	0.030	0.052	0.082
100000	0.319	0.441	0.760
1000000	2.362	3.506	5.945
100000000	24.26	46.93	79.01

- Time for collecting MIN, MAX, SUM, MEAN increases significantly from trials 4 to 5

- Time for assigning elements to respective buckets increases significantly during same trials

These are the two places where **FOR LOOPS** are incorporated; other elements in the program are simple, one-step functions

# POSSIBLE SOLUTIONS???

- NUMPY ARRAYS
  - HAVE EMBEDDED MIN, MAX, RANGE, MEAN FUNCTIONS WHICH CAN ELIMINATE THE FIRST FOR LOOP

```
29     intbucket1=[[0 for x in range(1000)] for y in range(buckets)]
30     numpy_start=datetime.datetime.now()
31     intBucketF=np.array(intBucket)
32     numpy_end=datetime.datetime.now()
33     stat_start=datetime.datetime.now()
34     Max=npamax(intBucketF)
35     Min=npamin(intBucketF)
36     Range=np.ptp(intBucketF)
37     Mean=np.mean(intBucketF)
```

# MORE DATA

Elements	Python Sort			New Sort - with NumPy Array		
	Start Time (s)	End Time (s)	Total Time (s)	Start Time (s)	End Time (s)	Total Time (s)
1000	57.98	57.98	0.00	4.02	4.02	0.00
10000	4.13	4.14	0.02	48.44	48.58	0.14
100000	25.69	25.78	0.10	13.70	14.93	1.23
1000000	15.06	18.03	2.98	0.59	15.80	16.40

- Run time has significantly decreased from before so the FOR LOOP was indeed affecting run time
- Can time still be cut using numpy array for the 2D list as well?

# MORE DATA

Python - Sort with Complete NumPy (Sort, Arrays)						
Elements	Trials	Time for NumPy Array (s)	Time for Stat Collection (s)	Time for Bucketing (s)	Total Time (s)	Sort
1000	AVG	0	0.000	0.023	0.023	0.000
10000	AVG	0	0.000	0.157	0.162	0.005
100000	AVG	0.031	0.000	1.661	1.926	0.265
1000000	AVG	0.267	0.013	16.530	45.109	28.580

- Time to create numpy array is **gradually increasing – may cause problems when data set increases** by factor of 100
- Time for **stat collection is negligible** – suggests that numpy array **does make it more efficient** in this aspect
- Time for bucketing with using numpy array
- Why don't the run times improve as much as expected?
  - **QUICKSORT** IS BEING USED TO SORT BUCKETS – **EMBEDDED PYTHON PROGRAM**

# IS PYTHON THE RIGHT LANGUAGE?

```
48     for x in range(upperlimit):
49         value=intBucketF[x] #calling value at index of list
50         if Mean<value: #determining which range it will fall into
51             bucketSelect=(int)(math.ceil(bucketSelectLower*(value-Min))) #seeing where
52         else:
53             bucketSelect=lowerGroups+(int)(math.ceil(bucketSelectUpper*(value-Mean)))
54             bucketSelect=(int)((value-Min)/Range)*upperlimit/1000) #reinstantiating bucket
55             bucketing_end=datetime.datetime.now()
56             for y in range(buckets):
57                 intbucket1.sort(intbucket1[y])
58             print(np.matrix(intbucket1))
```

- CHECKED TO SEE IF PROGRAM WAS ACTUALLY WORKING THE WAY IT WAS
  - TRIED MULTIPLE WAYS TO PRINT OUT THE 2D LIST AND THE FINAL SORTED ARRAY BUT UNSUCCESSFUL
    - MAYBE DUE TO FAULTY PROGRAMMING
    - DUE TO “BEGINNER’S LEVEL” IN THIS LANGUAGE

# DOES THIS PROGRAM EVEN WORK?

- IT DID WORK MAGNIFICENTLY IN JAVA (ADVANCED LEVEL PROGRAMMING SKILL IN LANGUAGE)
- DATA COLLECTED IN JAVA PROVES THAT IT DOES WORK
- PROOF OF CONCEPT IS ESSENTIAL TO MAKE SURE EFFORTS ARE NOT WASTED
- WILL JUST CONTINUE WITH COMPLETING FULL PROGRAM IN JAVA
  - **PROOF OF CONCEPT**

## NEXT STEPS

- APPLY MULTITHREADING – SORT BUCKETS ALL AT THE SAME TIME
- MODIFY ALGORITHM TO SORT VARIOUS TYPES OF DISTRIBUTIONS
  - BINOMIAL
  - POISSON
  - WEIBULL
  - LOG NORMAL
  - STUDENT'S T
  - HYPERGEOMETRIC
    - MERELY THE **FORMULA** NEEDS TO CHANGE
- CREATE A TOOL THAT WILL FIND DISTRIBUTION – IF NEEDED