2017

# An Evaluation of Sorting Algorithms

TCSS 558 – MASTER'S SEMINAR

JESSE BANNON

# Description of Sorting Algorithms

## Insertion Sort

Insertion sort constructs the sorted array by moving single elements right-to-left, comparing and moving itself left until it has satisfied its ordering condition. For example, suppose our list is [0, 2, 3, 1, 4]. [1] moves itself to the left until it's condition is satisfied:
[0, 2, 3, 1, 4]: (1 < 3) -> [0, 2, 1, 3, 4]: (1 < 2) -> [ 0, 1, 2, 3, 4]: (1 > 0)
The algorithm has a O(n^2) runtime complexity, and is computed in-place, which results in a O(n) space complexity.

## Selection Sort

Selection sort constructs the sorted array one element at a time starting on one side and progressing to the other side. For example, when sorting in ascending order from left-to-right, we start at index zero. We search every element in the list to find the min element in [0, n), and swap that element to reside in index zero. We repeat this process for [1, n), [2, n), …, [n-1, n) to sort the list. The algorithm has a O(n^2) runtime complexity, and is computed in-place, which results in a O(n) space complexity.

## Bubble Sort

Bubble sort naively iterates the entire list and swaps elements a single space per pass if the ordering condition does not satisfy a neighboring element until no swaps are performed in an iteration. The worst case is if the entire list is sorted except the last element should be the first element. To move the last element a single space, the algorithm must iterate the entire list. The algorithm has a O(n^2) runtime complexity, and is computed in-place, which results in a O(n) space complexity.

## Merge Sort

Merge sort is a deterministic algorithm where it partitions the data in a binary tree-like structure. A top-down merge sort recursively partitions the data in two until there are only one or two elements in each partition, then constructs a sorted partition by iterating the two sorted sub-partitions. This algorithm has a O(nlog(n)) runtime complexity, but requires an additional list of size n to construct sorted partitions, which results in a O(2n) = O(n) space complexity.

## Quick Sort

Quick sort is similar to merge sort in the sense that it recursively partitions the data in two. Instead of splitting equally in two, quick sort chooses a pivot element, and places elements less than the pivot element on the left and elements greater than the pivot element on the right. Once you have two partitions that are < pivot or > pivot, you recursively choose a pivot element for those new partitions and repeat the process until the list is sorted. On average, quicksort has a nlog(n) runtime complexity, but its worst case makes it O(n^2). Quick sort is computed in-place, which results in a O(n) space complexity.

# Description of Data Sets

*Shuffled (1000)* – elements [1,1000] shuffled
*Ascending (800)* – elements [1,800] in ascending (sorted) order
*Descending (800)* – elements [1,800] in descending order
*N-8 Equal Elements Shuffled (800)* – 9 unique elements in the data, where 8 of them have a count of 1
*Article Share Count* – Number of shares among 1000 articles published by Mashable
(https://archive.ics.uci.edu/ml/datasets/Online+News+Popularity)
*Breast Cancer Treatment Age* – Age of patients who received breast cancer treatment
(https://archive.ics.uci.edu/ml/datasets/Mammographic+Mass)


# Performance of Sorting Algorithms

Based on the performance metrics gathered below, there are several implications to consider when choosing a sorting algorithm. Insertion sort, selection sort, and bubble sort all fall into the class of naïve $O(n^2)$ sorting algorithms. Despite their high average complexity, they are still useful in certain situations. Because they operate in place and do not require recursion, they are especially useful for sorting small lists for several reasons: 1) They do not need to allocate additional address spaces/stacks because they do not recurse. 2) If they are contiguous in memory, fetching it is fast. 3) They can potentially be fast if the list is nearly sorted.

In the data where there are n-8 equal elements shuffled, insertion sort performed the best because it only had to consider those 8 elements and iteratively move them until their ordering condition was satisfied.

Both bubble sort and insertion sort performed well on the ascending (already ordered) data, because they iteratively check every element to see if any need to be swapped. If the list is sorted, the algorithms only have to perform n comparisons.

Selection sort is the worst of the three naive algorithms, because no matter what ordering of the data is, it must perform $O(n^2)$ comparisons before it completes. Hence why its performance line is relatively in the same position regardless of the data.

Merge sort is reliable in the sense that its runtime and space is deterministic for every run. Its low runtime complexity compared to other sorting algorithms makes it excel in performance no matter what data it sorts on. In nearly every case, merge sort had the best or was very close to the best algorithm in terms of runtime. Its only drawback is its memory complexity and recursive nature when allocating new address spaces/stacks.

Quick sort can easily vary in performance depending on the data and pivot policy. In my implementation, I chose the highest element in the partition as the pivot, which resulted in it being the slowest algorithm for the ascending data since every element landed in the left partition. In the uniform random scenarios, which is most often the case for real world data, it can excel in runtime without the burden of merge sort's space complexity.

# Issues with Metrics

For my measure of sortedness per algorithm, I used number of swaps or sets (merge sort). This metric however did not provide any useful information, because in most cases, performance is affected based on number of comparisons – selection sort can have 0 swaps but have $O(n^2)$ comparisons. It has no correlation with runtime or memory used, so I placed those graphs in a scatterplot since no line can be drawn to represent the metrics.

Memory usage was constant per algorithm based on the input data. I am not sure why this is something worth comparing/plotting for each data set when it could easily be represented better in a table.

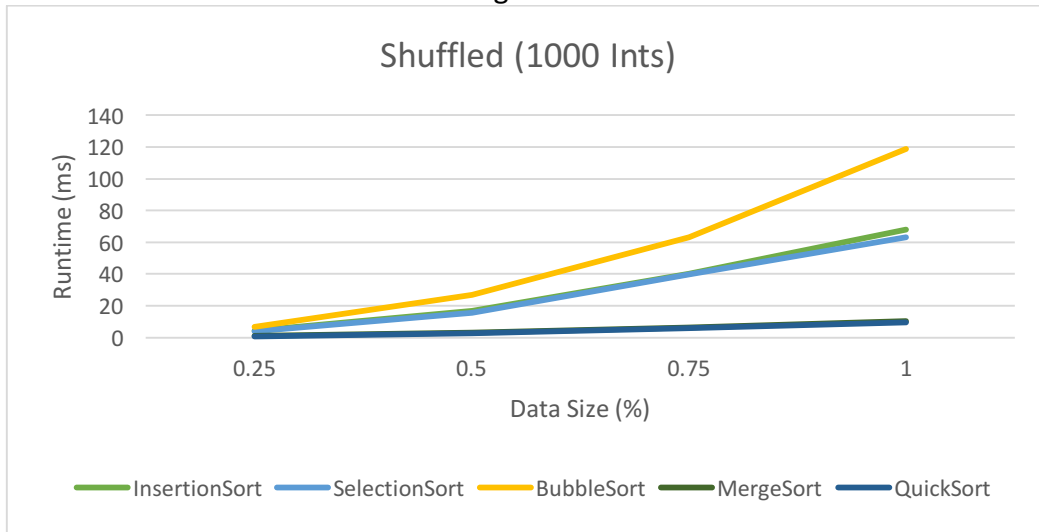The only useful performance metrics are Runtime against Data Size.

# Acknowledgements

# Metrics

## Runtime against Data Size

### Shuffled (1000 Ints)



Y-axis: Runtime (ms) — 0, 20, 40, 60, 80, 100, 120, 140
X-axis: Data Size (%) — 0.25, 0.5, 0.75, 1

Legend: InsertionSort, SelectionSort, BubbleSort, MergeSort, QuickSort

### Ascending (800 Ints)



Y-axis: Runtime (ms) — 0, 20, 40, 60, 80
X-axis: Data Size (%) — 0.25, 0.5, 0.75, 1

Legend: InsertionSort, SelectionSort, BubbleSort, MergeSort, QuickSort

### Descending (800 Ints)



Y-axis: Runtime (ms) — 0, 20, 40, 60, 80, 100, 120
X-axis: Data Size (%) — 0.25, 0.5, 0.75, 1

Legend: InsertionSort, SelectionSort, BubbleSort, MergeSort, QuickSort

# N - 8 Equal Elements Shuffled (800 Ints)



Runtime (ms) vs Data Size (%)

Legend: InsertionSort, SelectionSort, BubbleSort, MergeSort, QuickSort

# Article Share Count (1000 Ints)



Runtime (ms) vs Data Size (%)

Legend: InsertionSort, SelectionSort, BubbleSort, MergeSort, QuickSort

# Breast Cancer Treatment Age (961 Ints)



Runtime (ms) vs Data Size (%)

Legend: InsertionSort, SelectionSort, BubbleSort, MergeSort, QuickSort

# Runtime vs Swap/Set Count

## Shuffled (1000 Ints)

Runtime (ms) vs Swap/Set Count

- 118.6714
- 63.3136
- 67.9192
- 10.4824

## Ascending (800 Ints)

Runtime (ms) vs Swap/Set Count

- 42.7882
- 5.9708
- 4.8886

## Descending (800 Ints)

Runtime (ms) vs Swap/Set Count

- 112.5388
- 84.933
- 47.5572
- 5.5244

## N - 8 Equal Elements Shuffled (800 Ints)

Runtime (ms) vs Swap/Set Count

- 47.8302
- 43.8124
- 4.6788
- 5.0728

## Article Share Count (1000 Ints)

Runtime (ms) vs Swap/Set Count

- 117.7842
- 63.2786
- 62.1824
- 9.7544

## Breast Cancer Treatment Age (961 Ints)

Runtime (ms) vs Swap/Set Count

- 111.215
- 58.0386
- 58.607
- 9.274

# Memory Usage vs Data Size

## Shuffled (1000 Ints)



Memory (kb) vs Data Size (%)

InsertionSort — SelectionSort — BubbleSort — MergeSort — QuickSort

## Ascending (800 Ints)



Memory (kb) vs Data Size (%)

InsertionSort — SelectionSort — BubbleSort — MergeSort — QuickSort

## Descending (800 Ints)



Memory (kb) vs Data Size (%)

InsertionSort — SelectionSort — BubbleSort — MergeSort — QuickSort

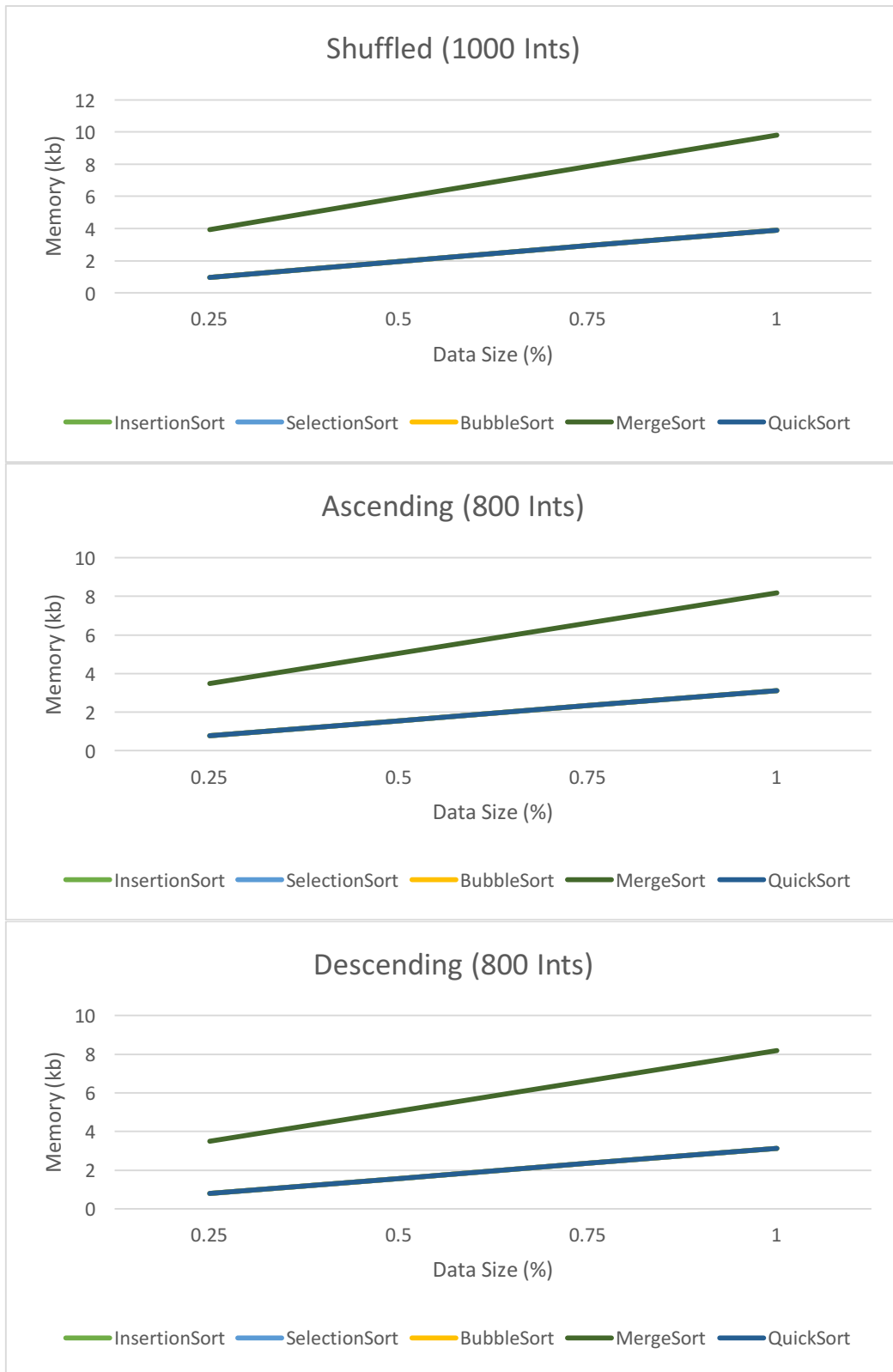# N - 8 Equal Elements Shuffled (800 Ints)



# Article Share Count (1000 Ints)
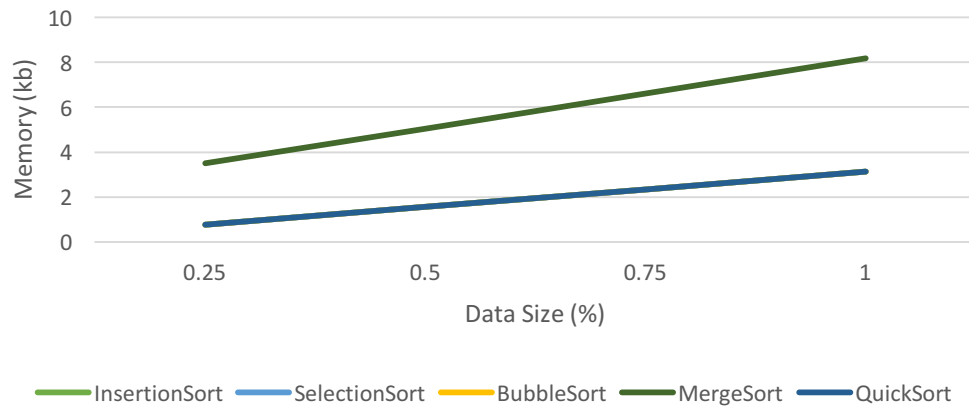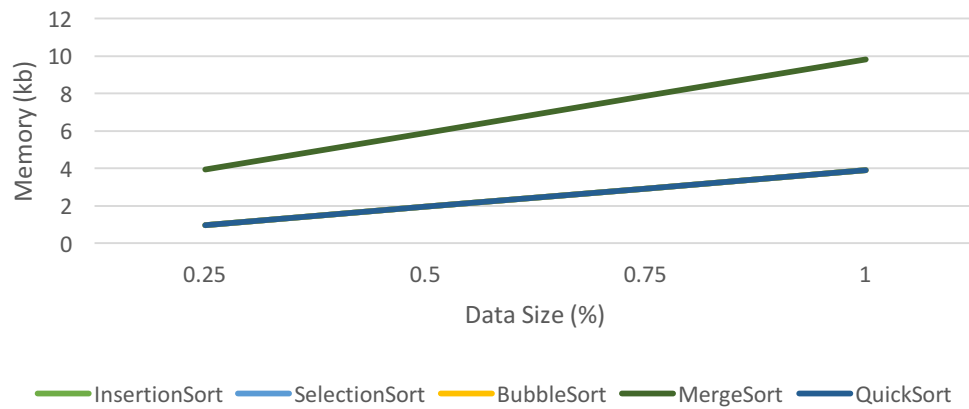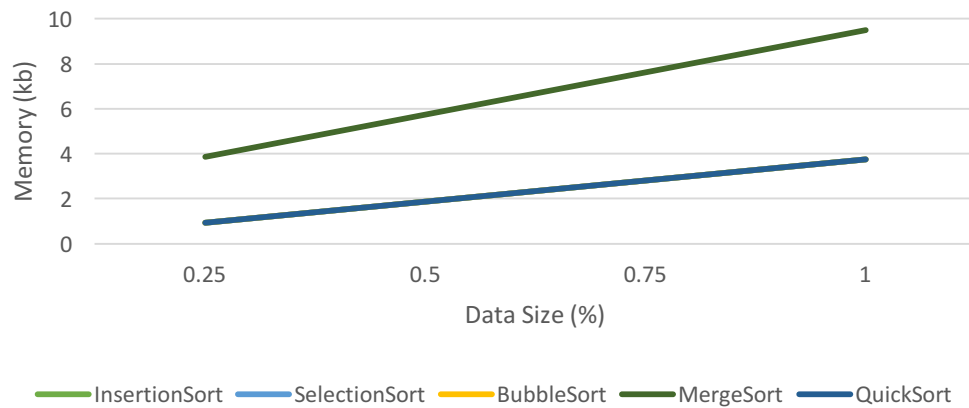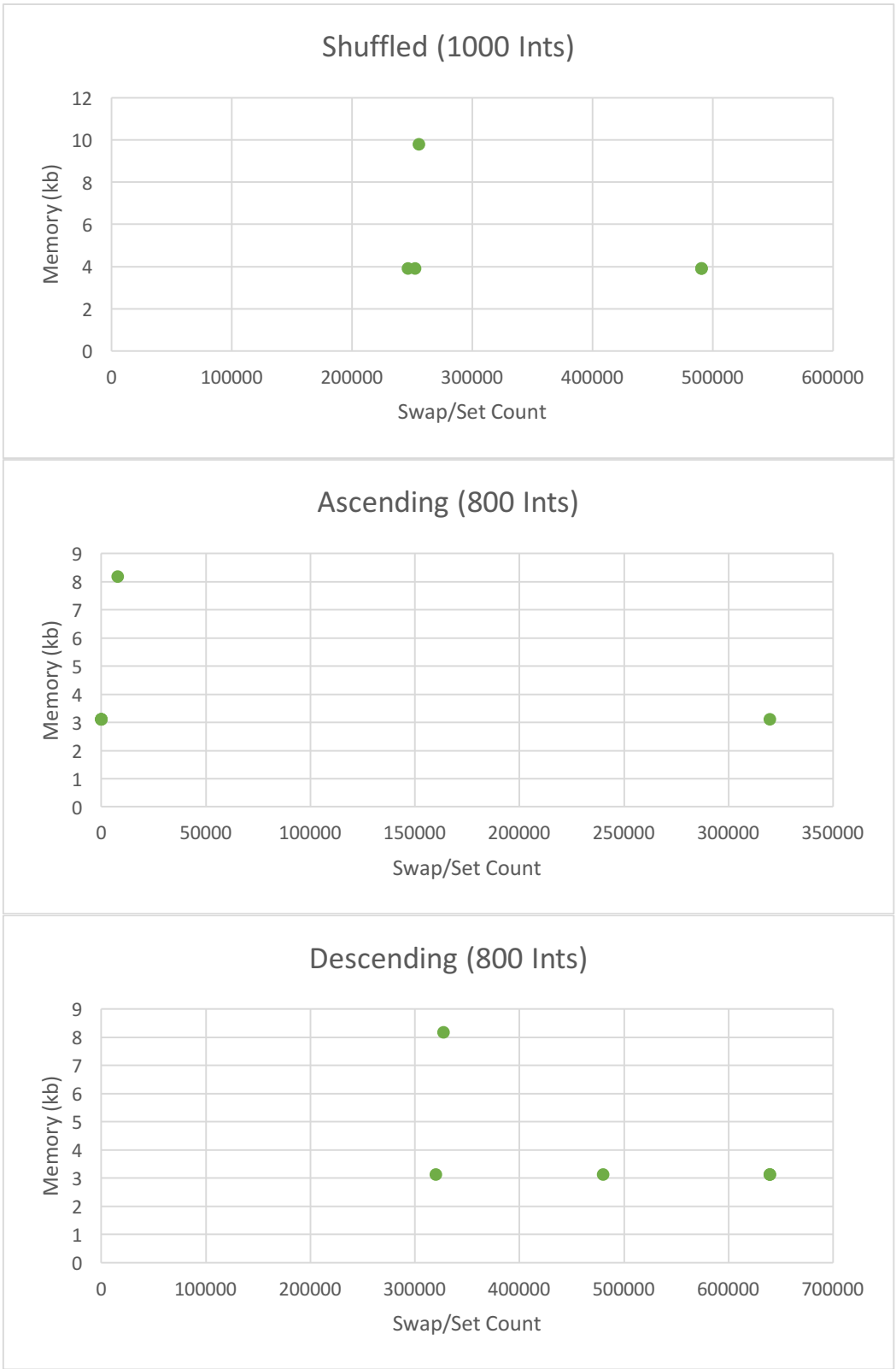


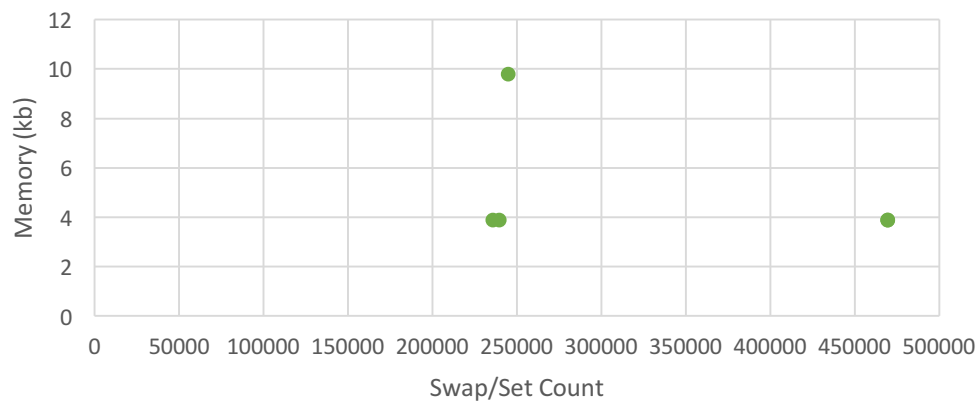# Breast Cancer Treatment Age (961 Ints)

# Memory Usage vs Swap/Set Count

## Shuffled (1000 Ints)



## Ascending (800 Ints)



## Descending (800 Ints)

## N - 8 Equal Elements Shuffled (800 Ints)

(Scatter plot: X-axis "Swap/Set Count" ranging 0 to 12000; Y-axis "Memory (kb)" ranging 0 to 9. Data points at approximately (4200, 3.1), (4600, 3.1), (6900, 3.1), and (11300, 8.2).)

## Article Share Count (1000 Ints)

(Scatter plot: X-axis "Swap/Set Count" ranging 0 to 500000; Y-axis "Memory (kb)" ranging 0 to 12. Data points at approximately (245000, 9.8), (240000, 3.9), and (470000, 3.9).)

## Breast Cancer Treatment Age (961 Ints)

(Scatter plot: X-axis "Swap/Set Count" ranging 0 to 500000; Y-axis "Memory (kb)" ranging 0 to 10. Data points at approximately (235000, 9.5), (230000, 3.8), and (450000, 3.8).)