

## Comp7505 Assignment 2

### Problem 2

With regard to my implementation of the “branch-and-bound” variation of `findMedianKmer`, the worst-case running time depends on two variables. With respect to the number of sequences, it is exponential where the exponent corresponds to the number of sequences (see Appendix 1.1). With respect to the length of the k-mer, it is quadratic (see Appendix 1.2).

Appendix 1.3 contains a chart illustrating the practical running times for 2-25 sequences when k ranges between 5 and 12. There is some missing data because the benchmark program would skip subsequent iterations if the previous took longer than 30 seconds.

It is possible to reach the conclusion of quadratic running time with respect to the length of the k-mer by analyzing the source code. The implementation gradually increases the size of the prefix until size k but for analysis we can assume that the size begins at k and gradually decreases until 0. This gives us the recurrence relation  $T(k) = O(S \cdot N \cdot k) + E \cdot T(k-1)$  where S is the number of DNA sequences, N is the length of each sequence, E is the length of the alphabet and k is the length of the k-mer (see Appendix 1.4). To simplify the analysis, assume that S, N, and E are held constant (in practice, E and S are initialized before running the algorithm so this assumption has merit, however, it is possible to have several values of N but we do have control over this before the algorithm is run). With the stated assumptions, the recurrence relation simplifies to  $T(k) = O(k) + T(k-1)$  which has the solution  $T(k) = O(k^2)$ .

### Problem 3

The method responsible for updating the trie is `putKmer` which has a theoretical running time of  $O(E \cdot k)$  where  $k$  is the length of the k-mer and  $E$  is the length of its alphabet (see Appendix 2.1). Since the k-mer's alphabet is constant,  $O(E \cdot k) = O(k)$  which means `putKmer` is linear with respect to the length of the k-mer.

The method which accesses the counts from the completed trie is `getCount` which has a theoretical running time of  $O(k)$  (see Appendix 2.2). It is possible that we reach the end of the path before finding the k-mer but the worst-case is actually finding the k-mer since that will lead to a longer path.

### Problem 4

The practical running time of `findMedianKmer` after improvements have been made is much faster (see Appendix 4.1). Using the trie to choose an optimal branching order helps minimally compared to reusing distances in this case. The reason is that the alphabet we're using is small so choosing a sub-optimal branch has little impact however recalculating distances for every recursive call becomes extremely expensive as the number of sequences and length of the k-mer grows.

Specifically, `getDistance` and `getDistances` are  $O(n \cdot k)$  where  $n$  is the length of the sequence and  $k$  is the length of the k-mer. The result must be evaluated over all sequences, so, in practice every distance calculation runs in  $O(n \cdot k \cdot N)$  where  $N$  is the number of sequences. By reusing distance calculations, in my implementation, we only use `getDistances` in the above way once. Every other call to `getDistances` thereafter happens within `getExtDists` which always only calculates the distance for a k-mer of length 1 so the running time of `getDistances` within `getExtDists` is  $O(n \cdot N)$ . The method `getExtDists` then sums the new, suffix distance with the previously calculated prefix distance which is also  $O(n \cdot N)$ . Therefore, by reusing distances, we cut the running time of finding distances from  $O(n \cdot k \cdot N)$  to  $O(n \cdot N)$  which is a dramatic improvement.

## Appendix 1.1

Plots of running time (in milliseconds) vs. k-mer length where n is the number of sequences.

Figure 1 -- Plot for n = 1

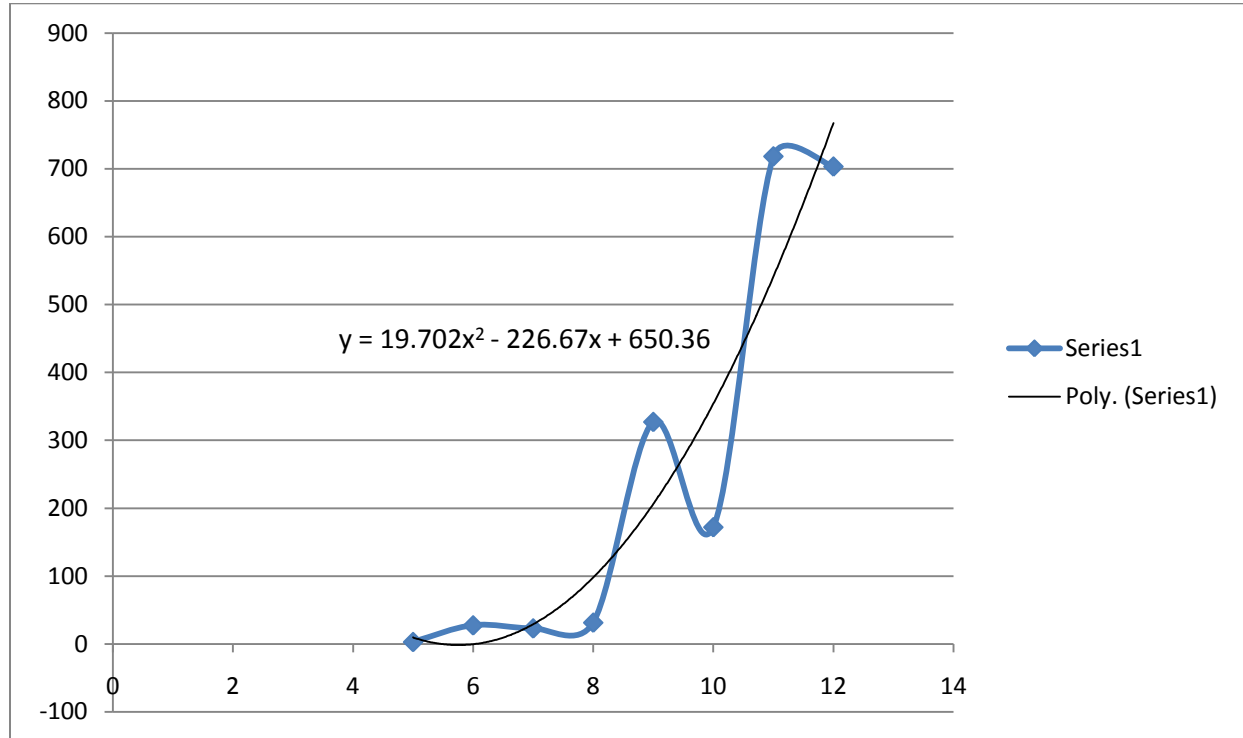


Figure 2 -- Plot for n = 2

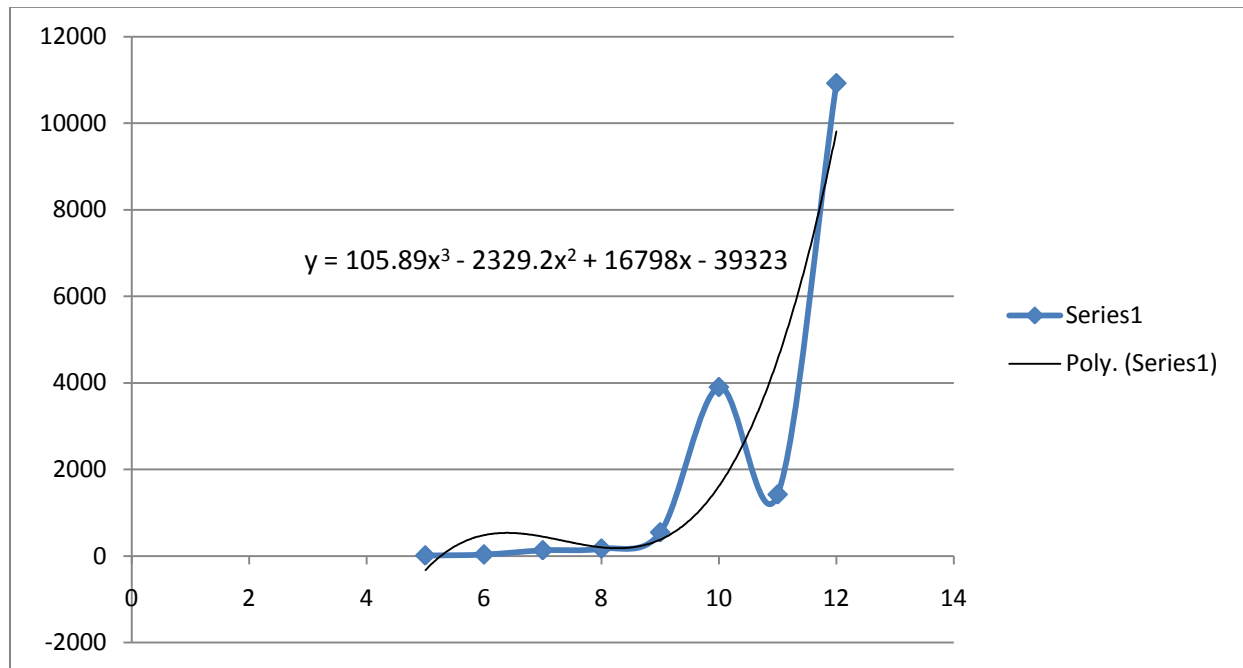


Figure 3 -- Plot for n = 4

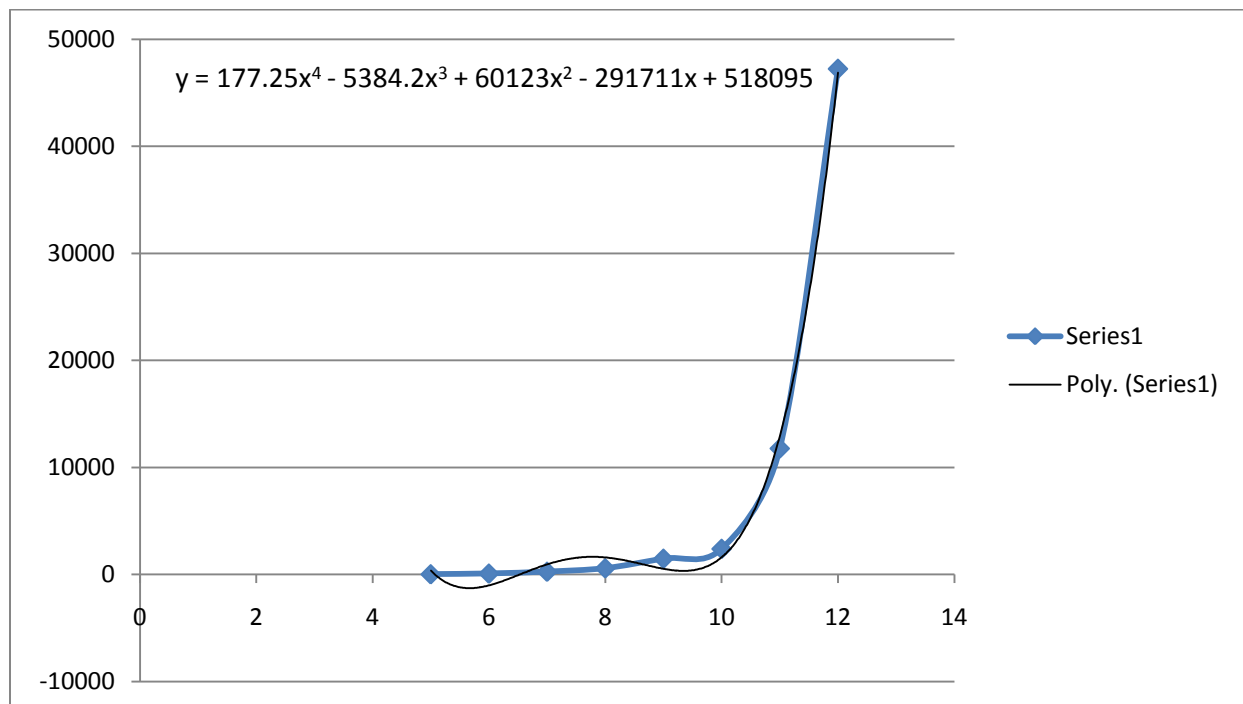


Figure 4 -- Plot for n = 5

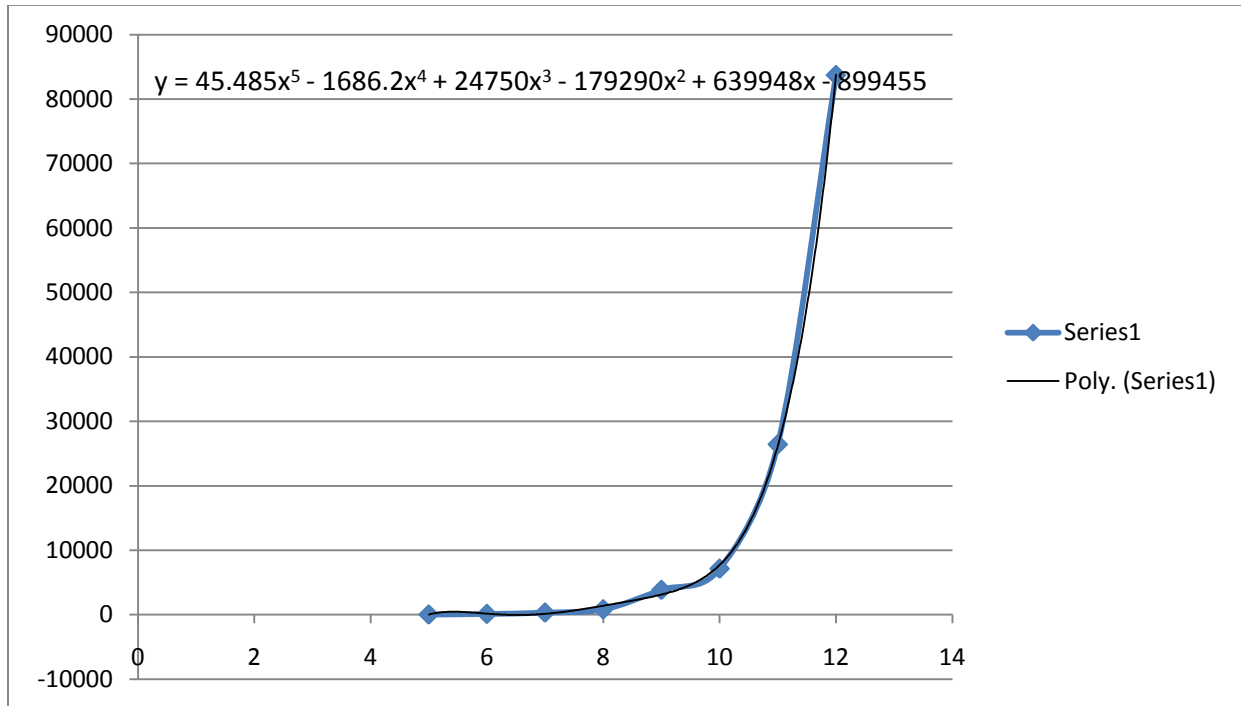
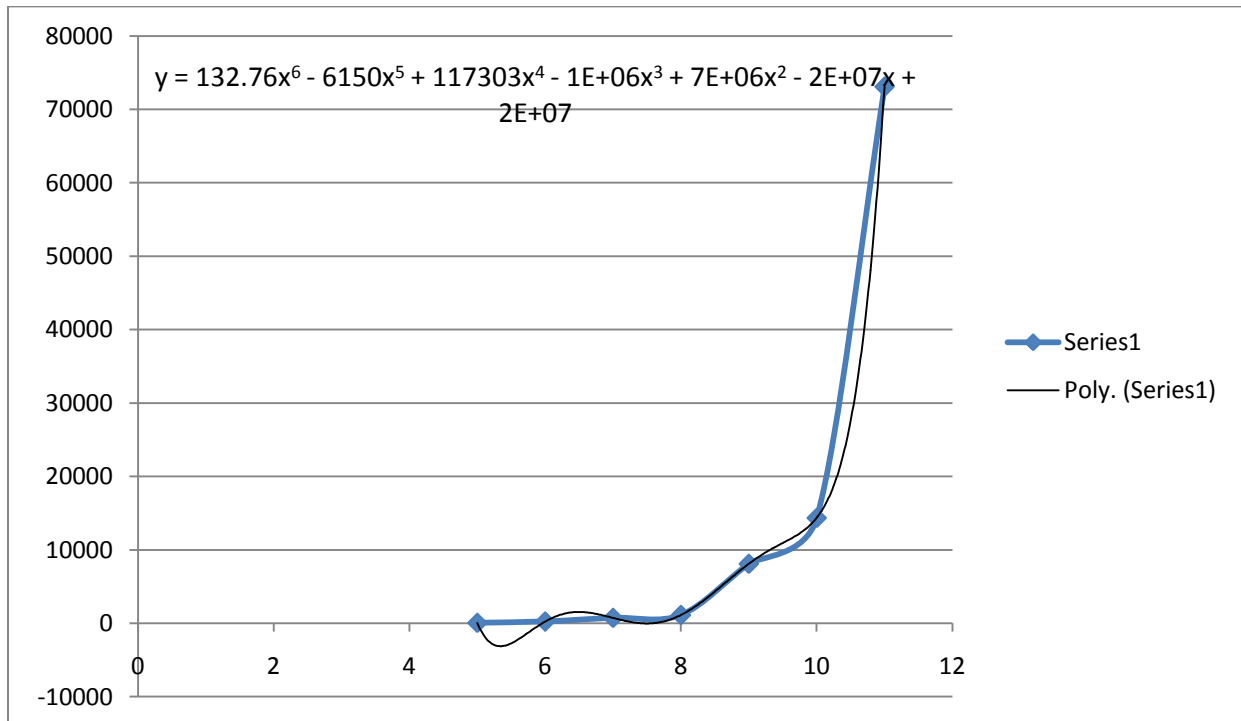


Figure 5 -- Plot for n = 6



## Appendix 1.2

Plots of running time (in milliseconds) vs. number of sequences where k is the length of the k-mer.

Figure 6 -- Plot for k = 5

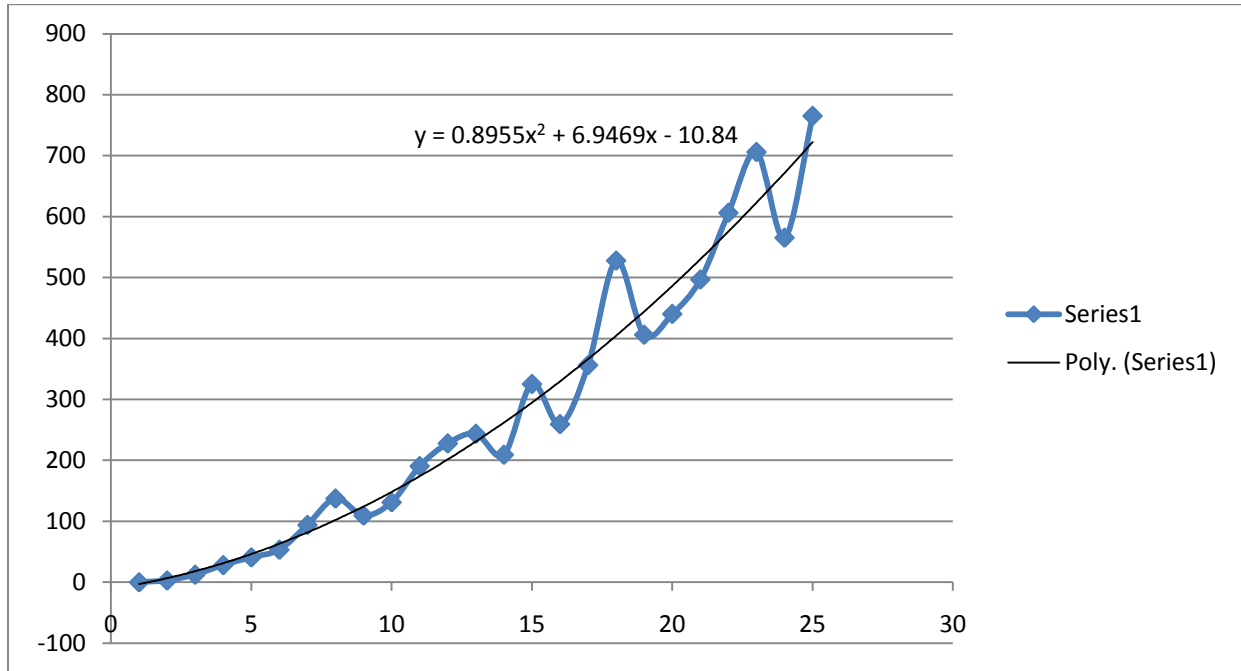


Figure 7 -- Plot for k = 6

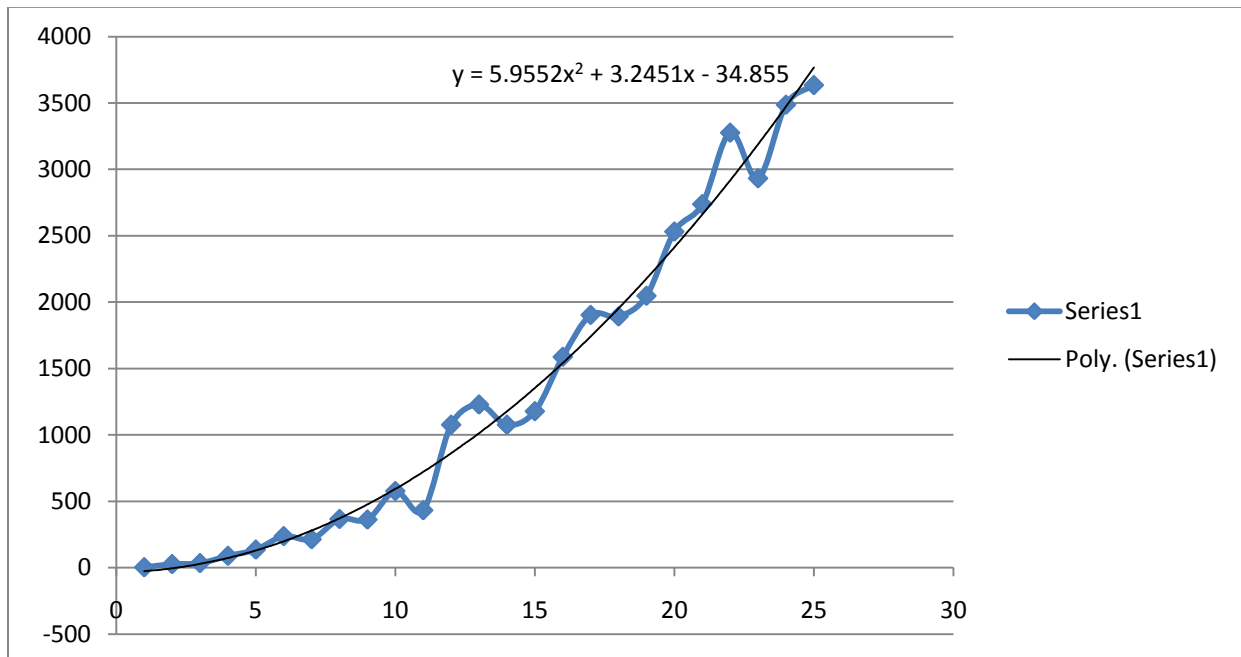


Figure 8 -- Plot for k = 7

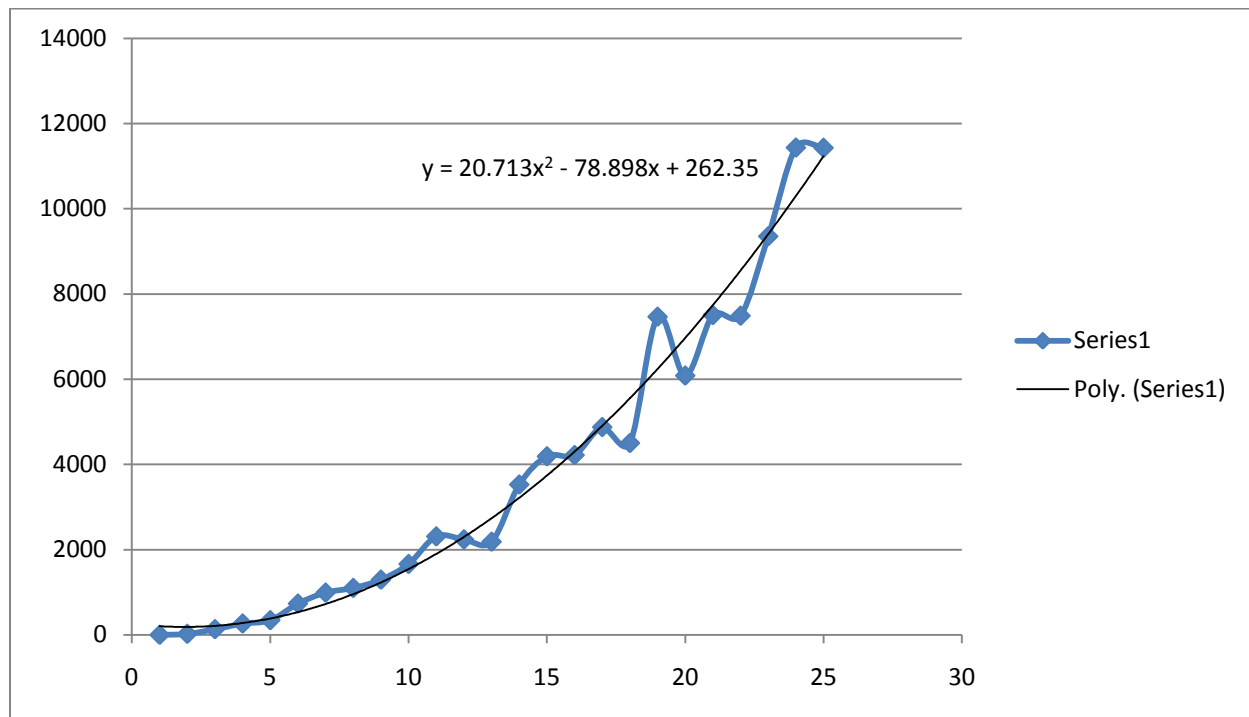
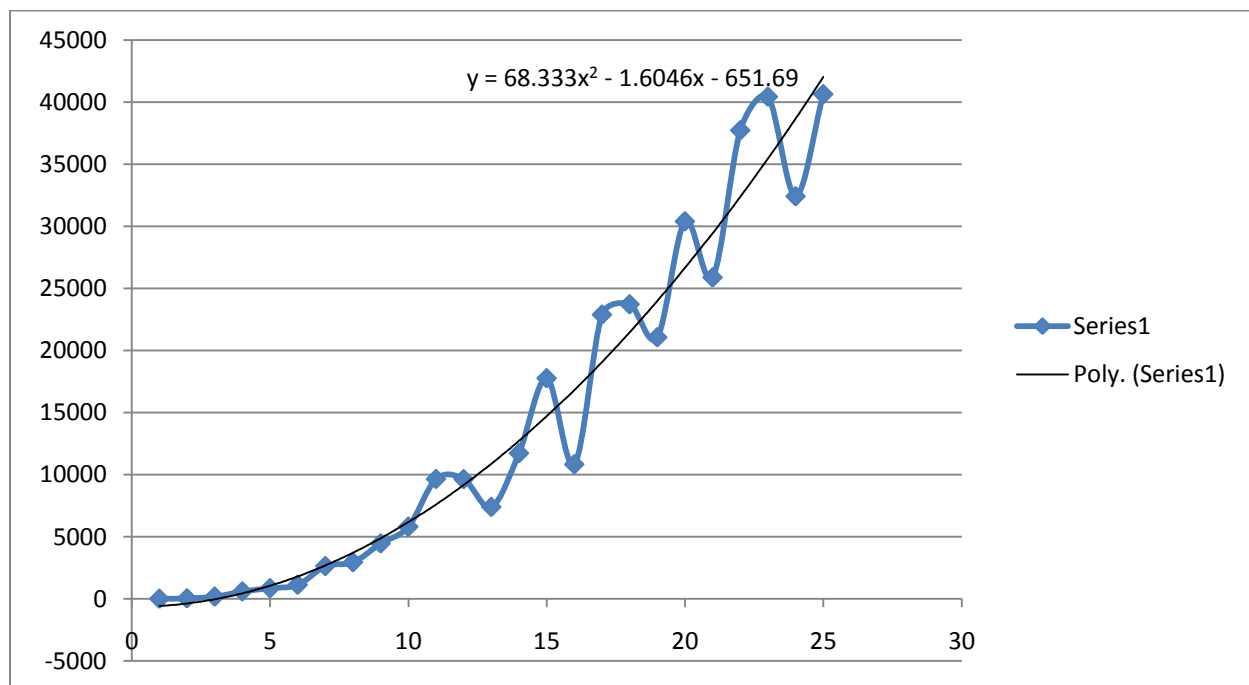
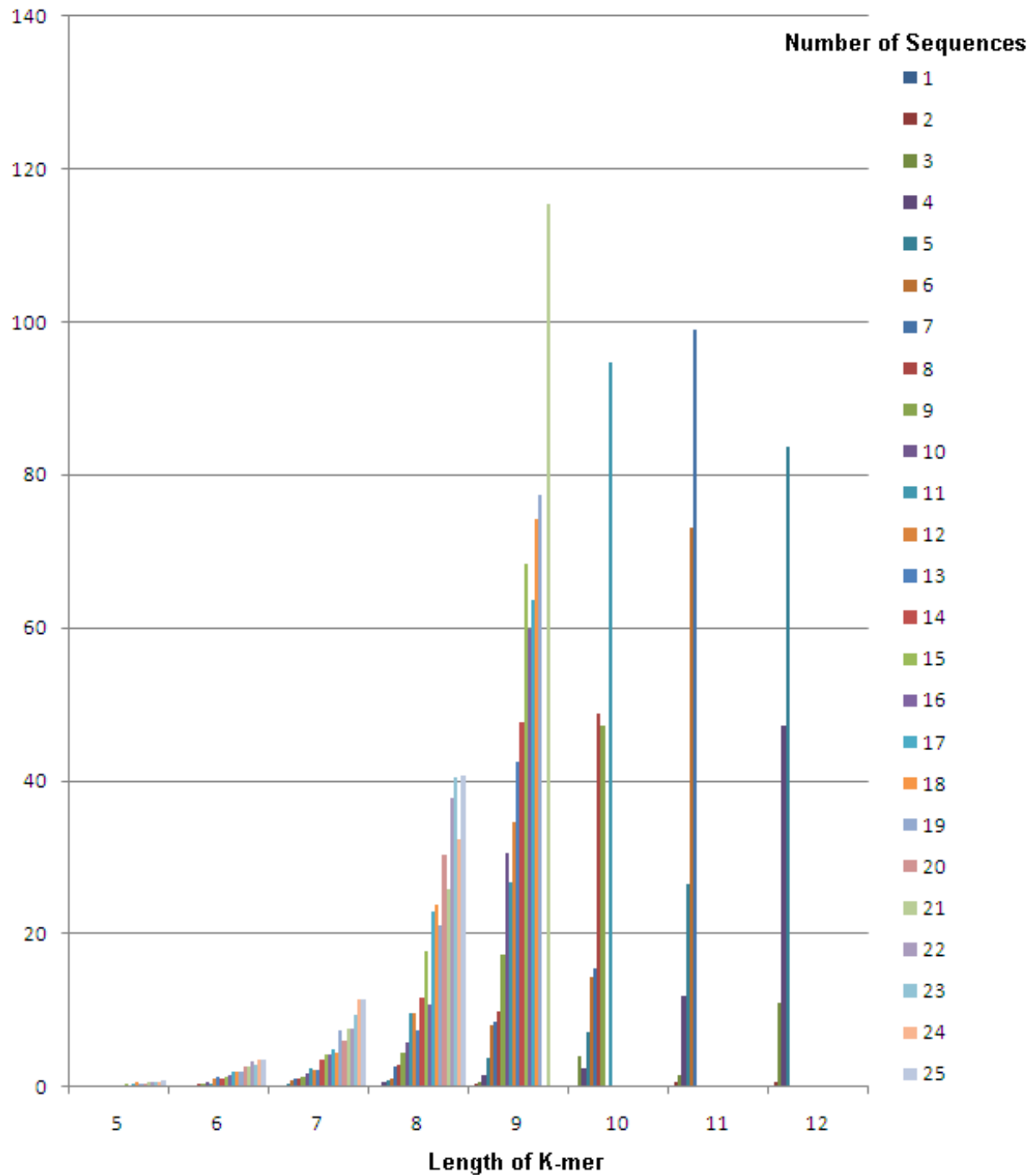


Figure 9 -- Plot for k = 8







## Appendix 1.4

Asymptotic analysis of findMedianKmer.

<code>public Distance findMedianKMer(KMer prefix, int k, Distance old_best)</code>	$T(k) = O(S*N*k) + E*T(k-1)$
<code>{</code>	
<code>    if (prefix == null)</code>	$O(1)$
<code>        prefix = new KMer(alpha, k);</code>	$O(k)$
<code>    if (prefix.isComplete()) {</code>	$O(1)$
<code>        int distance = 0;</code>	$O(1)$
<code>        for (DNASequence seq : this.seqs) {</code>	$O(S) * O(N*k) = O(S*N*k)$
<code>            distance += getDistance(seq, prefix);</code>	$O(N*k)$ where $N =  seq $ , $k =  prefix $
<code>        }</code>	
<code>        return new Distance(distance, prefix);</code>	$O(1)$
<code>    }</code>	
<code>    Distance best_dist = new Distance(k * this.seqs.length + 1, prefix);</code>	$O(1)$
<code>    if( old_best != null) {</code>	$O(1)$
<code>        int distance = 0;</code>	$O(1)$
<code>        for (DNASequence sequence : this.seqs)</code>	$O(S) * O(N*k) = O(S*N*k)$
<code>            distance += getDistance(sequence, prefix);</code>	$O(N*k)$ where $N =  seq $ , $k =  prefix $
<code>        if (distance &gt;= old_best.actual) {</code>	$O(1)$
<code>            return old_best;</code>	$O(1)$
<code>        }</code>	
<code>    }</code>	
<code>    Distance temp_dist;</code>	$O(1)$
<code>    for (KMer kmer : prefix.getExtensions()) {</code>	$O(E*?)$ Where $E =  \alpha $
<code>        /* Keep old_best throughout traversal */</code>	
<code>        if( old_best!=null) {</code>	$O(1)$
<code>            if (best_dist.actual &gt;= old_best.actual)</code>	$O(1)$
<code>                best_dist = old_best;</code>	$O(1)$
<code>            }</code>	
<code>        temp_dist = findMedianKMer(kmer, k, best_dist);</code>	$?$
<code>        if (temp_dist.actual &lt;= best_dist.actual) {</code>	$O(1)$
<code>            best_dist = temp_dist;</code>	$O(1)$
<code>        }</code>	
<code>    }</code>	
<code>    return best_dist;</code>	$O(1)$

## Appendix 2.1

Asymptotic analysis of putKmer.

---

<code>public void putKMer(TrieNode parent, KMer kmer) {</code>	
<code>    int[] symbols = kmer.getKMer();</code>	$O(1)$
<code>    for (int i = 0; i &lt; symbols.length; i++) {</code>	$O(E*k)$ where $k =  kmer $
<code>        TrieNode current = parent.children[symbols[i]];</code>	$O(1)$
<code>        if (current == null) {</code>	$O(1)$
<code>            current = new TrieNode(this.alpha);</code>	$O(E)$ where $E =  \text{alpha} $
<code>            current.symbol = symbols[i];</code>	$O(1)$
<code>            parent.children[symbols[i]] = current;</code>	$O(1)$
<code>        }</code>	
<code>        current.count++;</code>	$O(1)$
<code>        parent = current;</code>	$O(1)$
<code>    }</code>	
<code>}</code>	

## Appendix 2.2

Asymptotic analysis of getCount.

<code>public int getCount(KMer kmer) {</code>	$O(k)$
<code>/* The path to follow */</code>	
<code>int[] path = kmer.getKMer();</code>	$O(1)$
<code>/* The current node of the trie */</code>	
<code>TrieNode current = this.root;</code>	$O(1)$
<code>/*</code>	
<code> * The count. If there is nothing in the tree then it's not possible to</code>	
<code> * find the kmer so default to 0</code>	
<code> */</code>	
<code>int count = 0;</code>	$O(1)$
<code>/* The depth which has been traveled */</code>	
<code>int depth = 0;</code>	$O(1)$
<code>/*</code>	
<code> * Continue traversing the path until either we have traversed the whole</code>	
<code> * path or we have reached a leaf node</code>	
<code> */</code>	
<code>while (depth &lt; kmer.getLevel() &amp;&amp; current.children[path[depth]] != null) {</code>	$O(k)$
<code>    /* go to the correct child and increase the depth */</code>	
<code>    current = current.children[path[depth++]];</code>	$O(1)$
<code>}</code>	
<code>if (depth &lt; kmer.getLevel()) {</code>	$O(1)$
<code>    /* Followed the path but didn't find kmer */</code>	
<code>    return 0;</code>	$O(1)$
<code>} else {</code>	
<code>    /* Found kmer ending at current node */</code>	
<code>    count = current.count;</code>	$O(1)$
<code>}</code>	
<code>return count;</code>	$O(1)$
<code>}</code>	

## Appendix 3.1

Runtime (in seconds) vs Length of K-mer and Number of Sequences for implementation of `findMedianKmer` which leverages the trie data structure and re-using distances of parent nodes.

