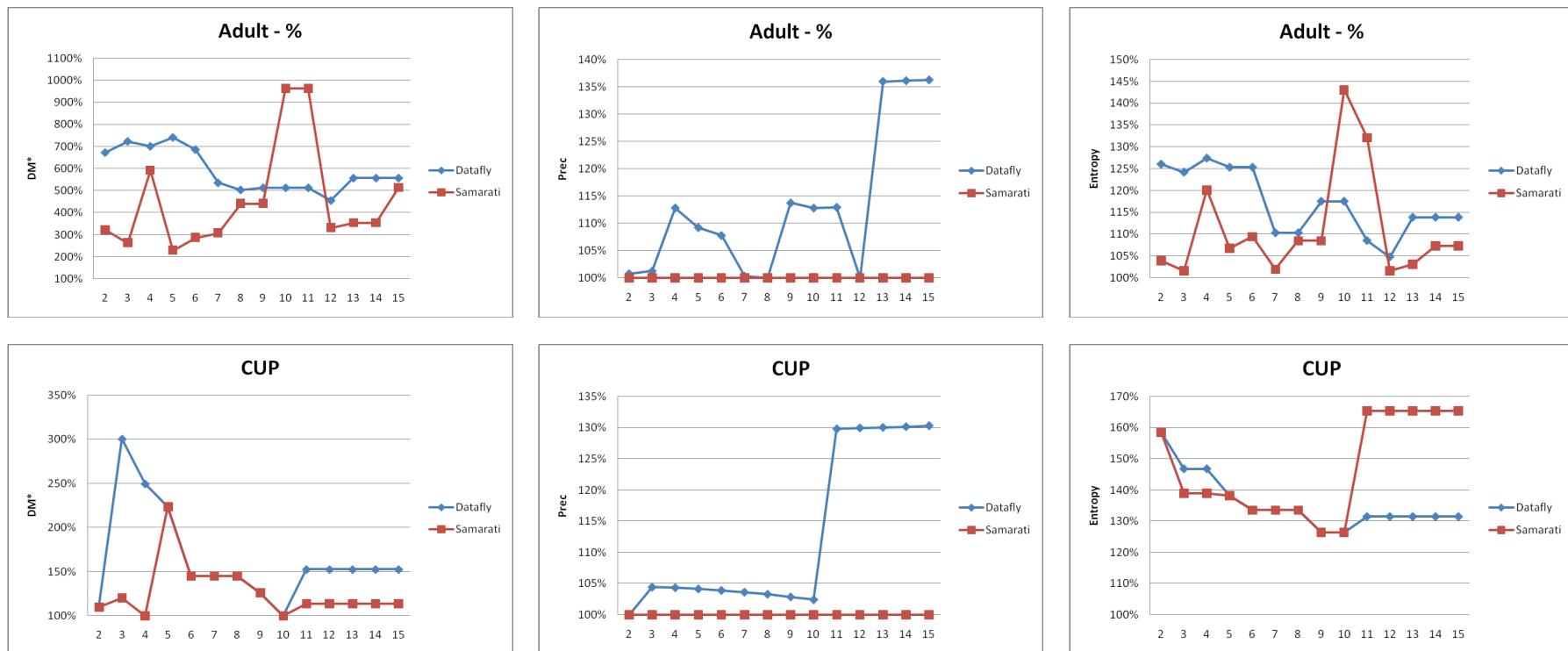


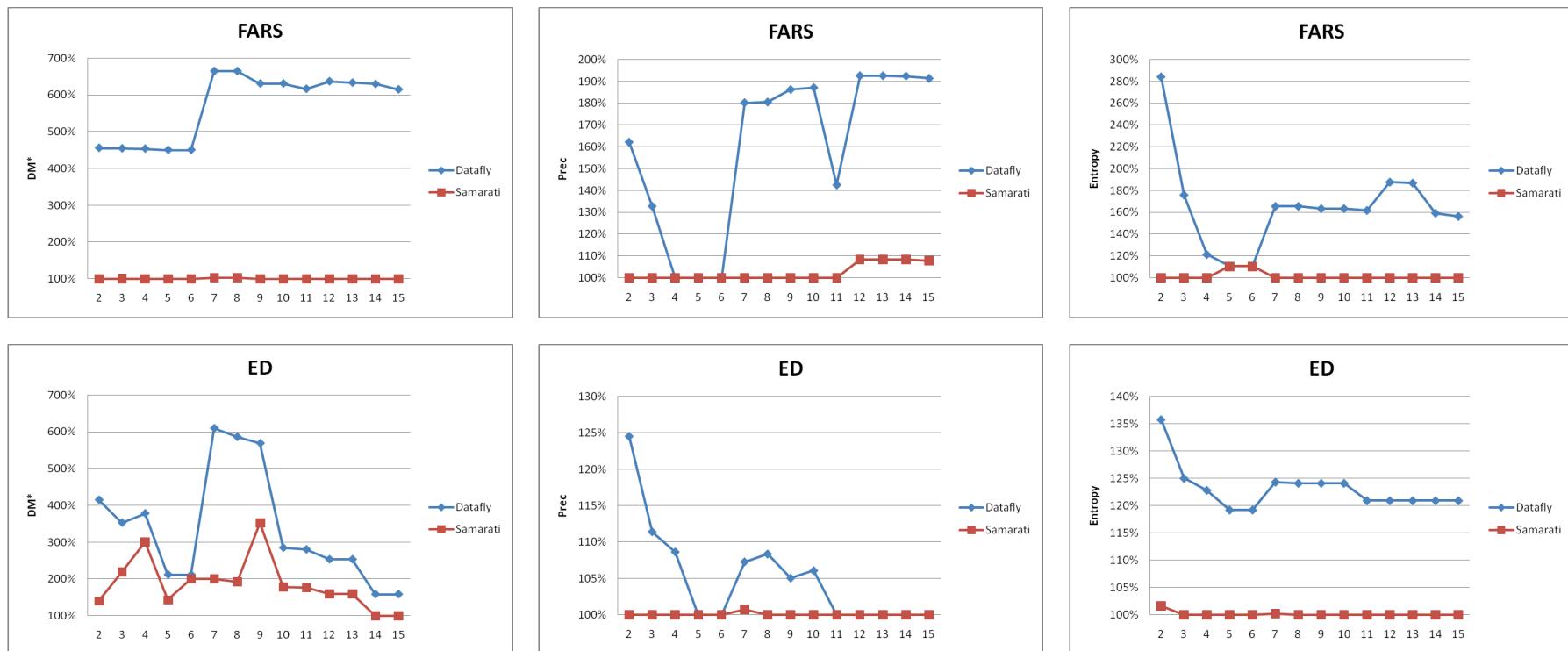
## **Appendix C, Part 1: On-line Tables and Figures for El Emam et al., “A Globally Optimal k-Anonymity Method for the De-identification of Health Data”**

Description	Quasi-identifiers	No. Records	No. of Nodes
<b>Adult</b> The adult dataset from the UC Irvine machine learning data repository. This is an extract from the US census: <a href="ftp://ftp.ics.uci.edu/pub/machine-learning-databases/adult">ftp://ftp.ics.uci.edu/pub/machine-learning-databases/adult</a>	<ul style="list-style-type: none"> <li>• Age (3)</li> <li>• Profession (2)</li> <li>• Education (2)</li> <li>• Marital status (2)</li> <li>• Position (2)</li> <li>• Race (1)</li> <li>• Sex (1)</li> <li>• Country (3)</li> </ul>	30,162	5,184
<b>FARS</b> Department of Transportation Fatal crash information: <a href="http://www-fars.ntsa.dot.gov/main.cfm">http://www-fars.ntsa.dot.gov/main.cfm</a>	<ul style="list-style-type: none"> <li>• Age (4)</li> <li>• Race (1)</li> <li>• Month of Death (3)</li> <li>• Day of Death (2)</li> </ul>	101,034	120
<b>CUP</b> Data from the Paralyzed Veterans Association on veterans with spinal cord injuries or disease: <a href="http://kdd.ics.uci.edu/databases/kddcup98/kddcup98.html">http://kdd.ics.uci.edu/databases/kddcup98/kddcup98.html</a>	<ul style="list-style-type: none"> <li>• ZIP code (5)</li> <li>• Age (4)</li> <li>• Gender (2)</li> <li>• Income (3)</li> </ul>	63,441	360
<b>Pharm</b> Prescription records from the Children's Hospital of Eastern Ontario pharmacy for 18 months. This is for inpatients only and excludes acute cases. A de-identified version of this data is disclosed to commercial data aggregators.	<ul style="list-style-type: none"> <li>• Age (4)</li> <li>• Postal code (FSA) (3)</li> <li>• Admission date (6)</li> <li>• Discharge date (6)</li> <li>• Sex (1)</li> </ul>	7,318	2,352
<b>ED</b> Emergency department records from Children's Hospital of Eastern Ontario for July 2008. This data is disclosed for the purpose of disease outbreak surveillance and bio-terrorism surveillance.	<ul style="list-style-type: none"> <li>• Admission date (3)</li> <li>• Admission time (7)</li> <li>• Postal Code (7)</li> <li>• Date of Birth (6)</li> <li>• Sex (1)</li> </ul>	4,090	3,584
<b>Niday</b> A registry of all newborns in Ontario for 2006-2007: <a href="https://www.nidaydatabase.com/info/index.shtml">https://www.nidaydatabase.com/info/index.shtml</a>	<ul style="list-style-type: none"> <li>• Maternal postal code (7)</li> <li>• Baby DoB (4)</li> <li>• Mother DoB (6)</li> <li>• Baby sex (2)</li> <li>• Aboriginal status (2)</li> <li>• Primary language (2)</li> </ul>	124,933	7,560

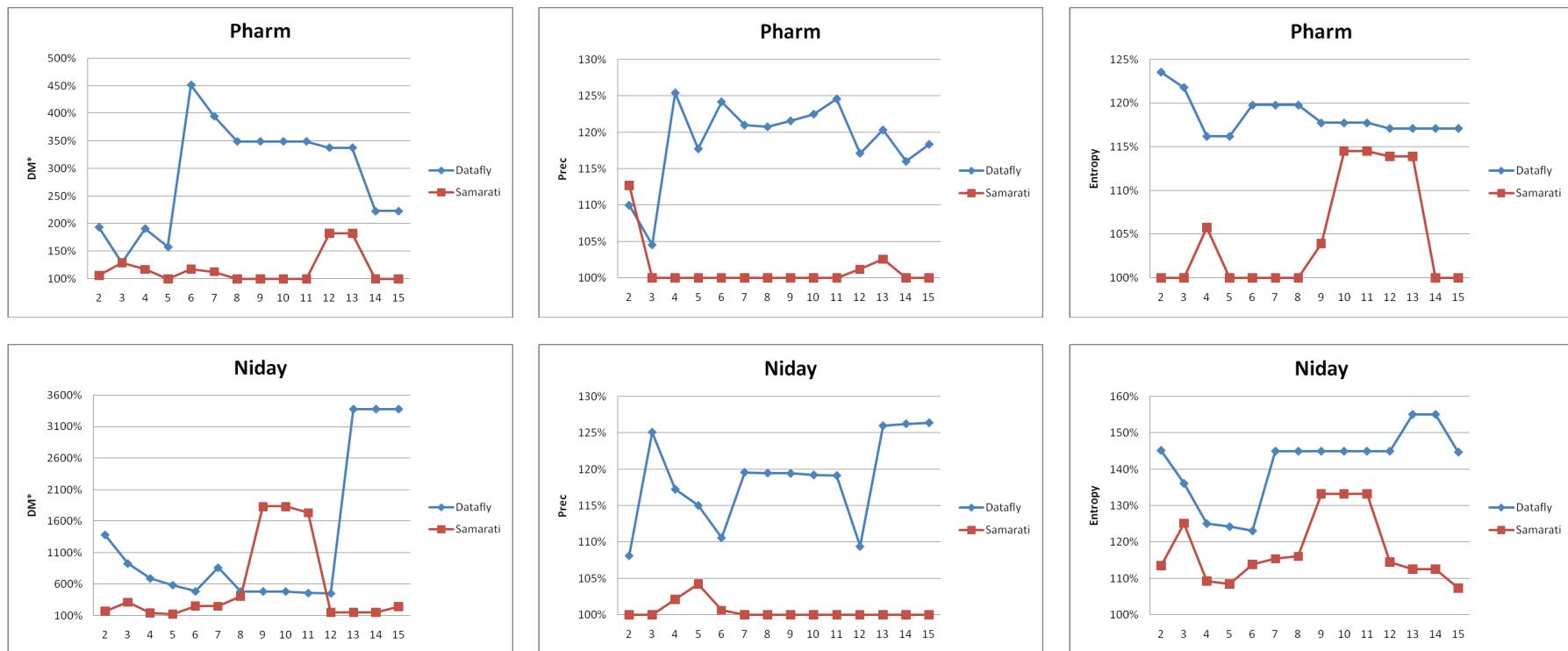
**Table 1:** Summary of the six datasets that we used in our comparative evaluation. For each dataset the quasi-identifiers used and the height of the generalization hierarchy are shown in the second column, the number of observations in the original dataset in column 3, and the number of nodes in the overall lattice in the last column.



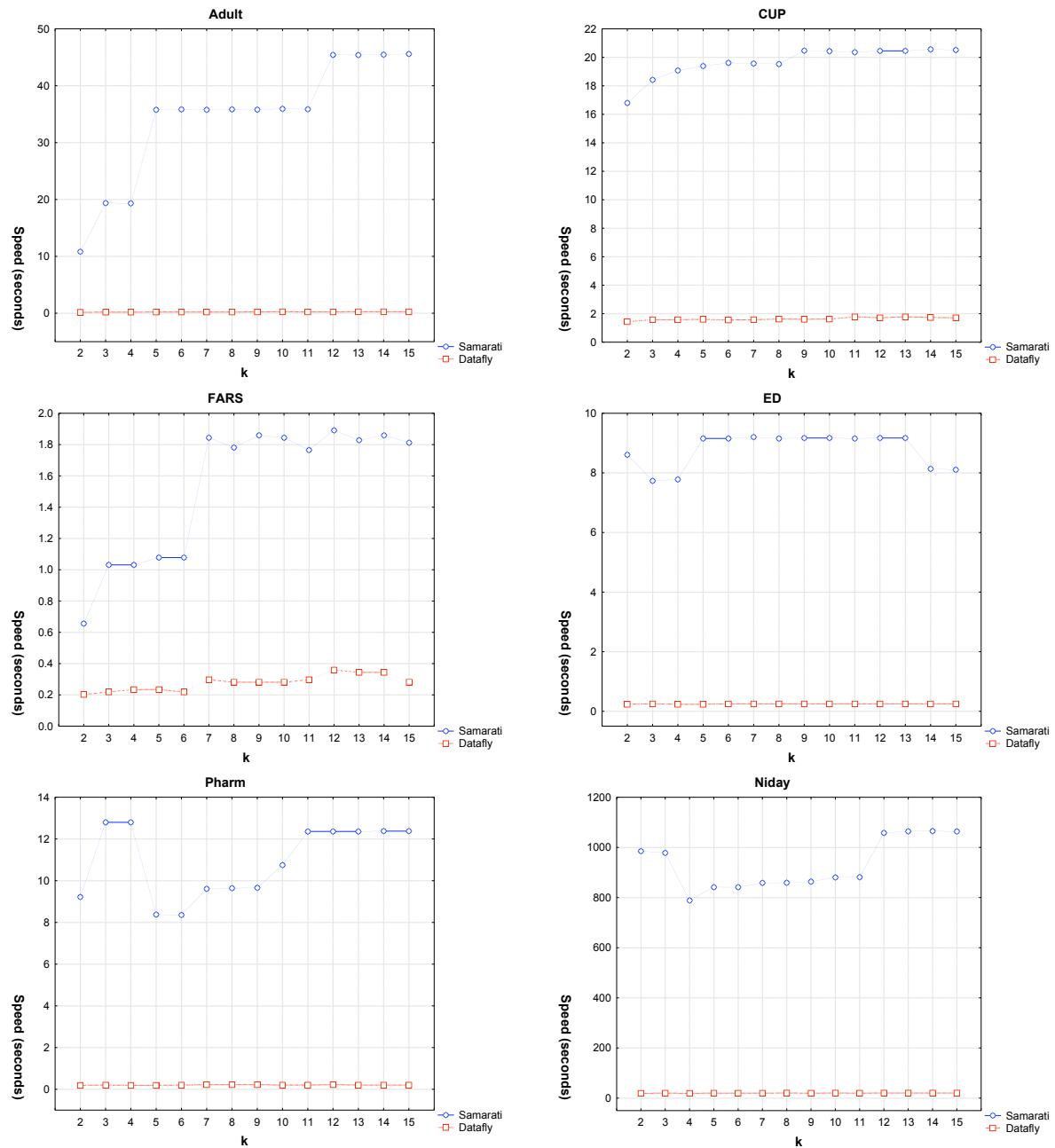
**Figure 6a:** Information loss results for the six datasets. The x-axis shows the value of  $k$  as it was varied from 2 to 15. The y-axis shows the information loss relative to our algorithm. The graphs show the results for a 5% suppression limit.



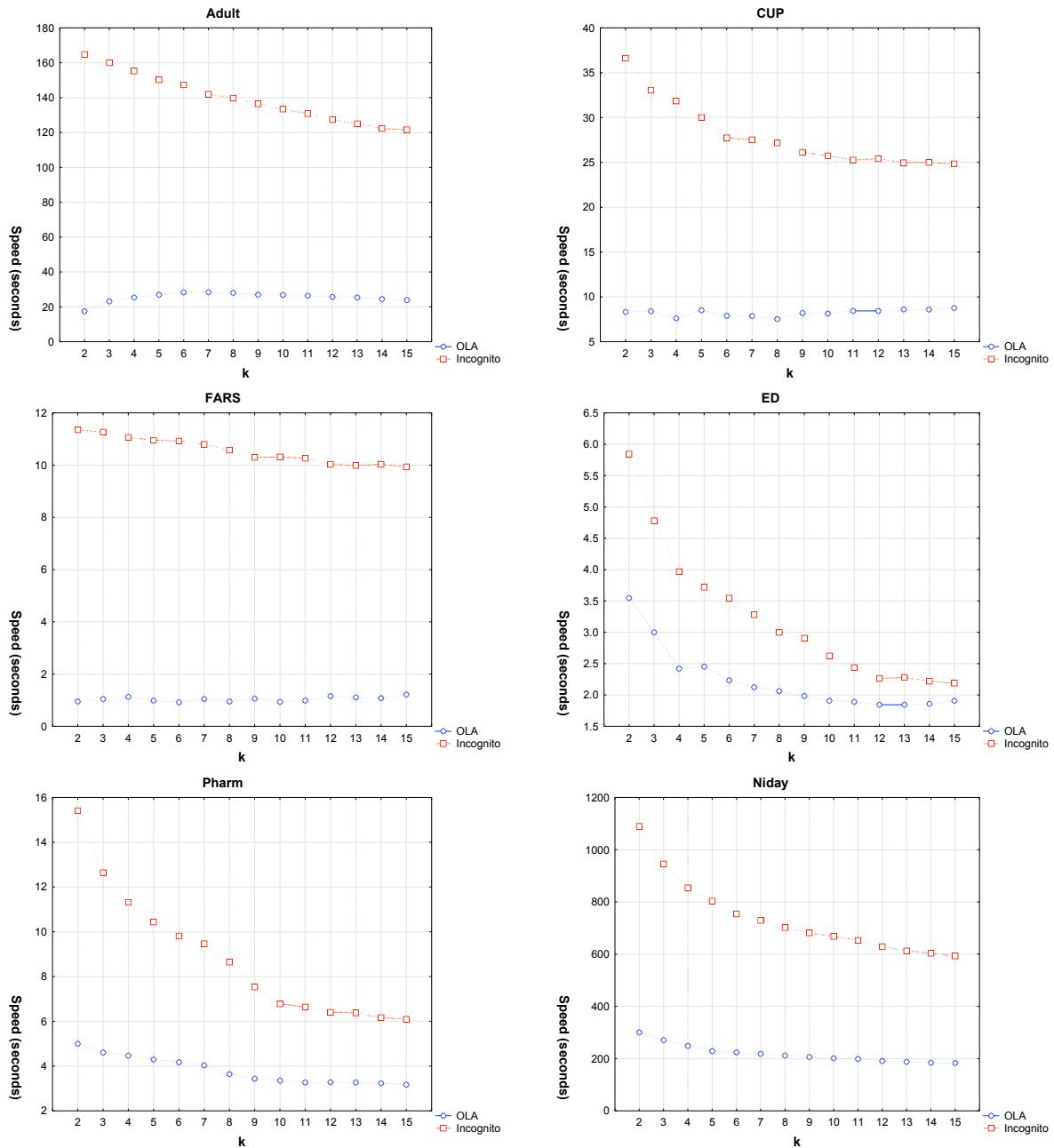
**Figure 6b:** Information loss results for the six datasets. The x-axis shows the value of  $k$  as it was varied from 2 to 15. The y-axis shows the information loss relative to our algorithm. The graphs show the results for a 5% suppression limit.



**Figure 6c:** Information loss results for the six datasets. The x-axis shows the value of  $k$  as it was varied from 2 to 15. The y-axis shows the information loss relative to our algorithm. The graphs show the results for a 5% suppression limit.



**Figure 7:** The speed to obtain a solution for Datafly and Samarati on the six datasets across values of  $k$  2 to 15 and a maximum suppression value of 5%. The timing is in seconds and includes all pre- and post- processing that may be needed to run each algorithm. These results were obtained on a dual core 3.2 GHz Pentium D processor with 3GB of RAM.

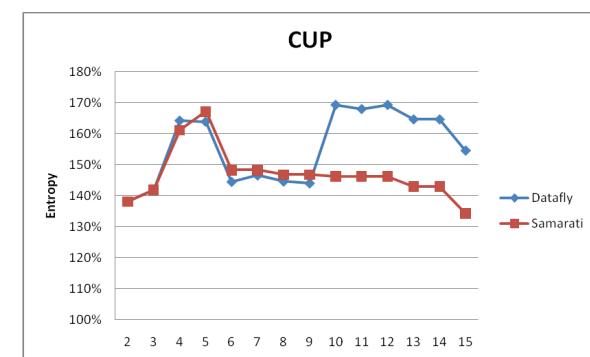
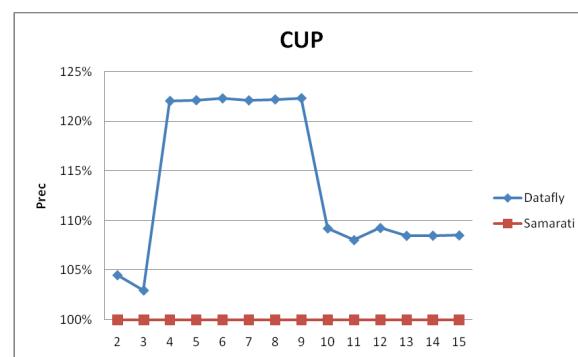
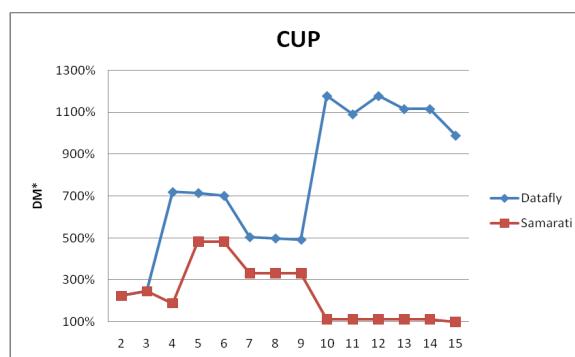
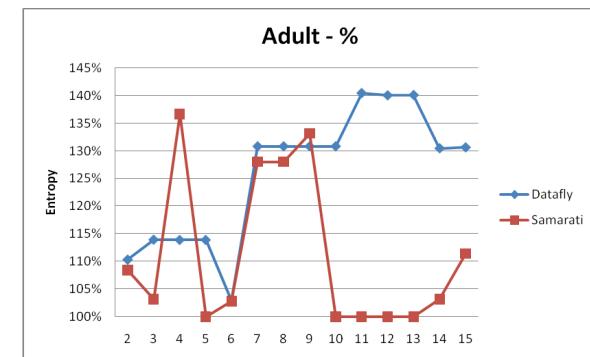
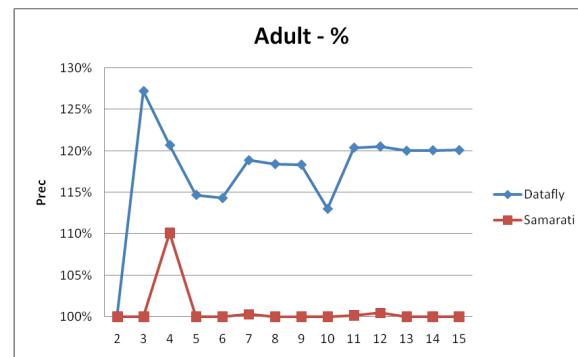
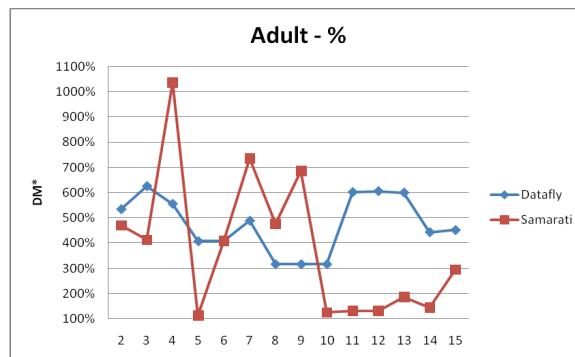


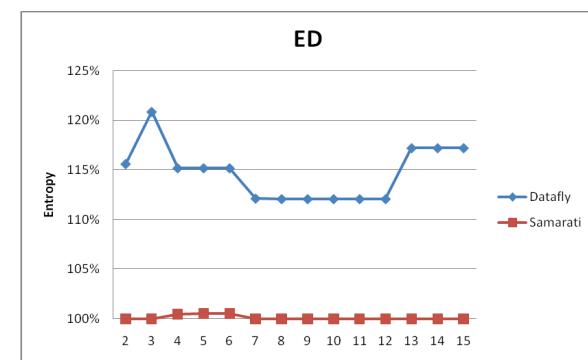
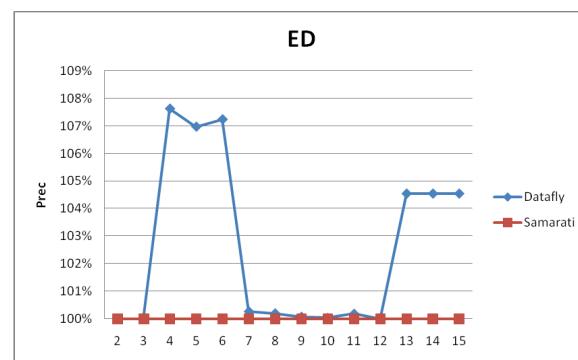
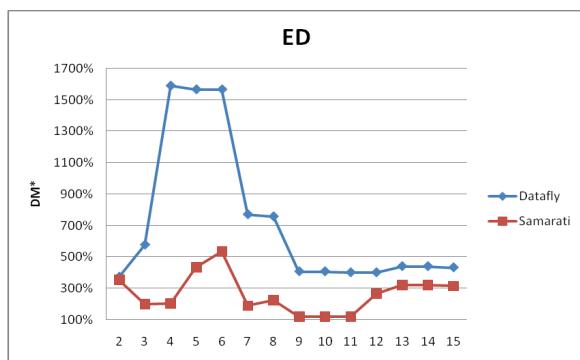
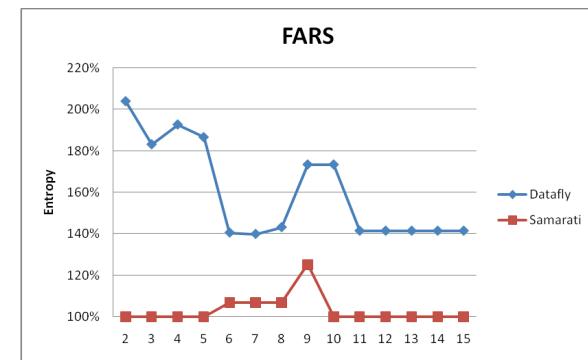
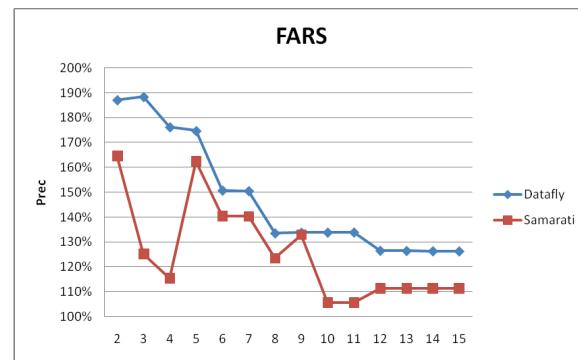
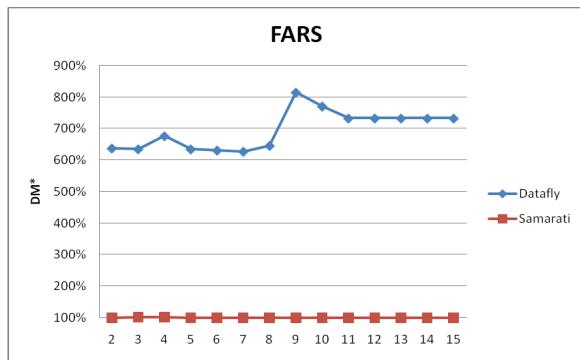
**Figure 9:** The speed to obtain a globally optimal solution using OLA and Incognito on the six datasets across values of  $k$  2 to 15 and a maximum suppression value of 5%. The timing is in seconds and includes all pre- and post-processing that may be needed to run each algorithm. These results were obtained on a dual core 3.2 GHz Pentium D processor with 3GB of RAM.

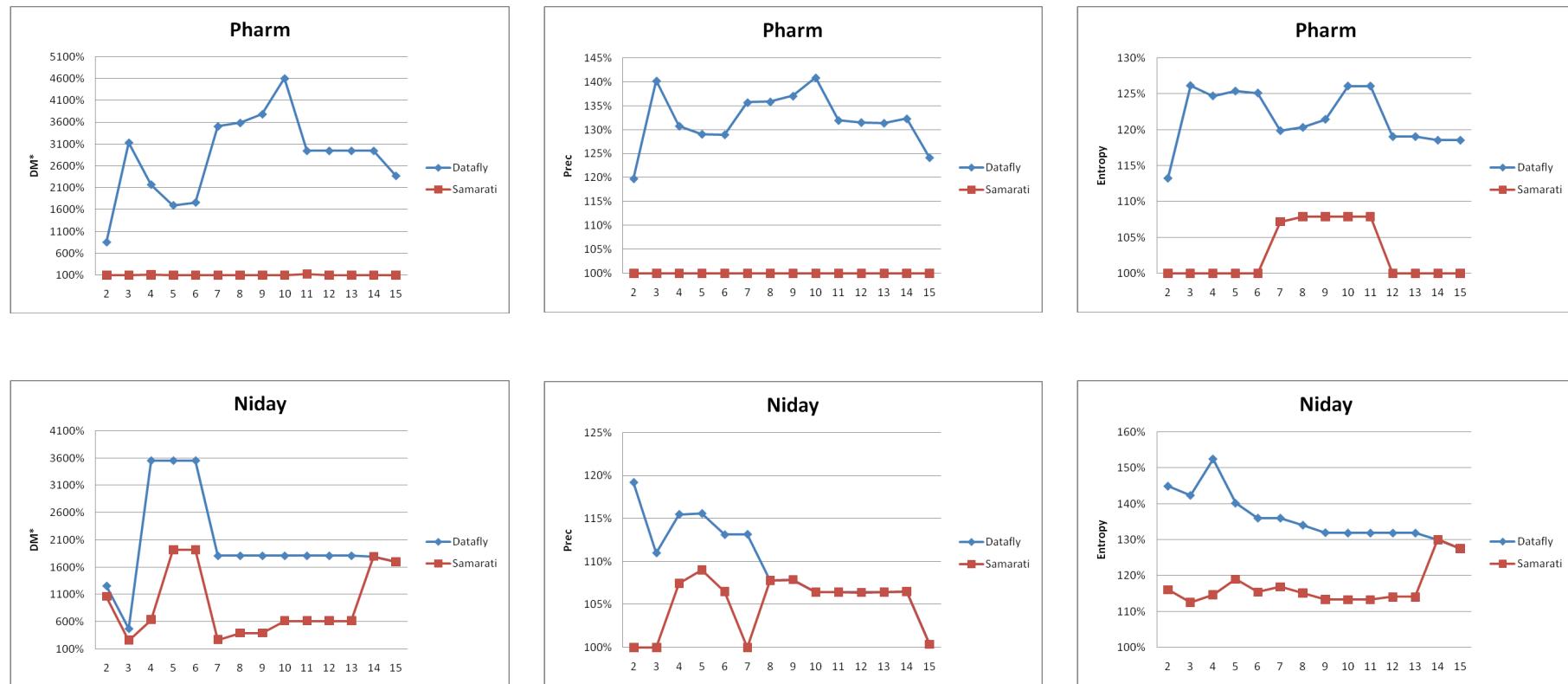
## Appendix C, Part 2: Results for 1% and 10% Suppression Limits

In this section we present all of the results graphs for the 1% and 10% suppression limits. The way these graphs should be interpreted is described in the main body of the paper. These results are consistent with the 5% suppression limit results.

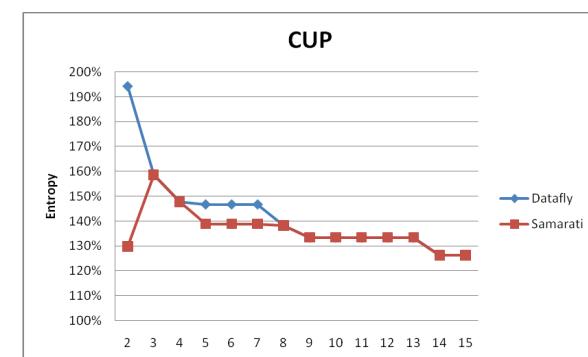
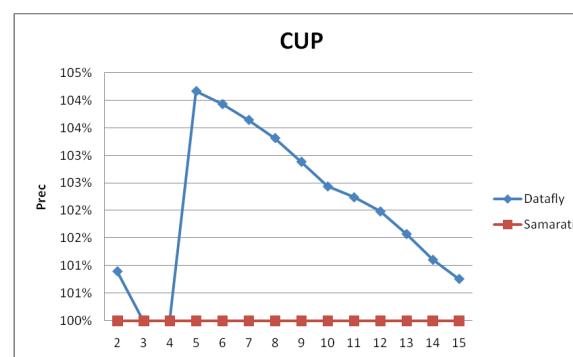
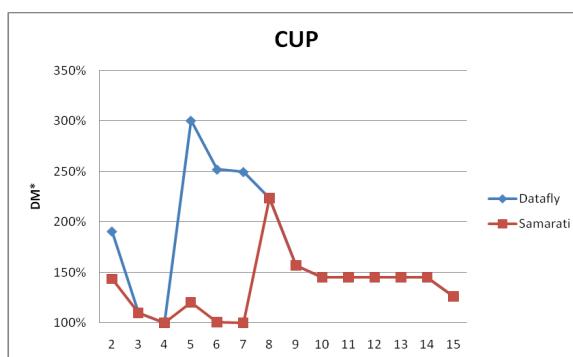
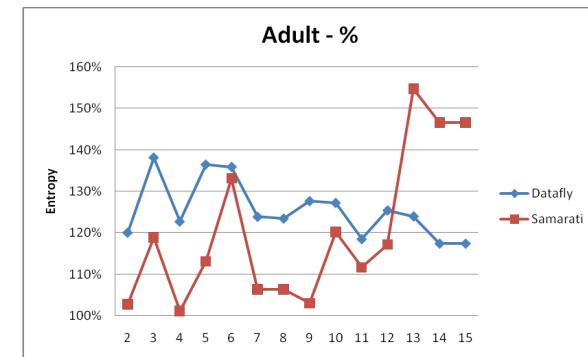
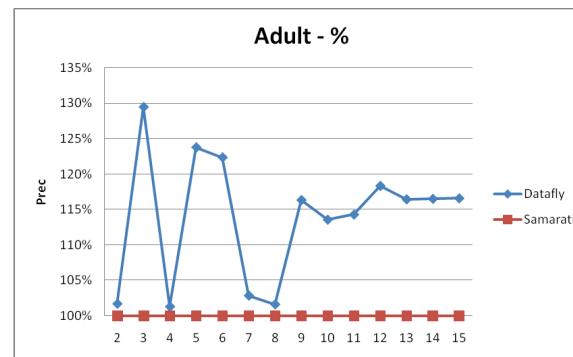
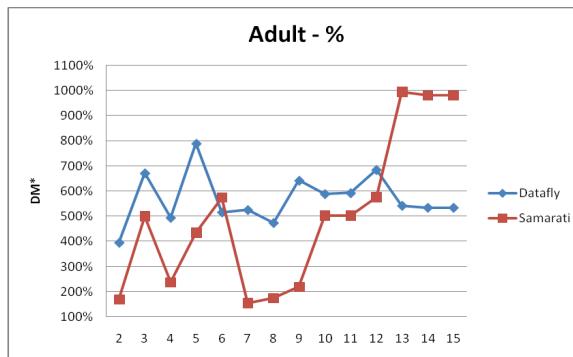
### Information Loss Results for the 1% Suppression Limit

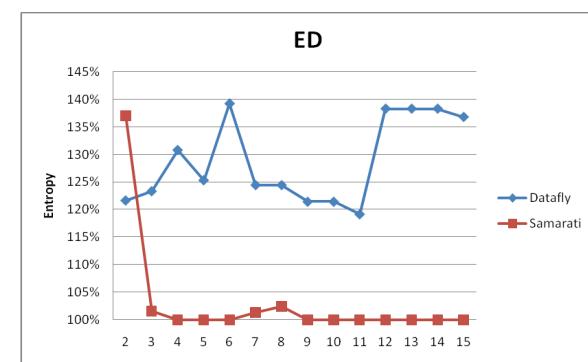
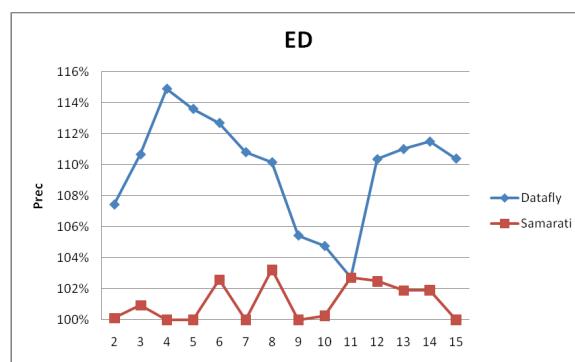
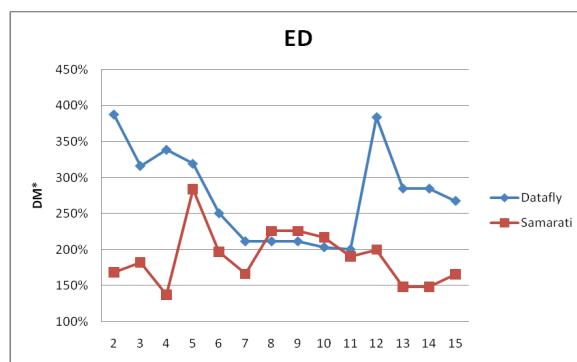
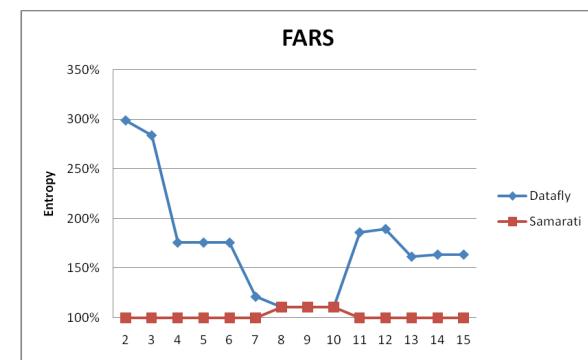
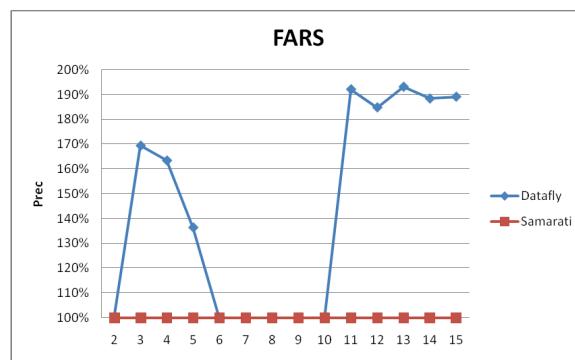
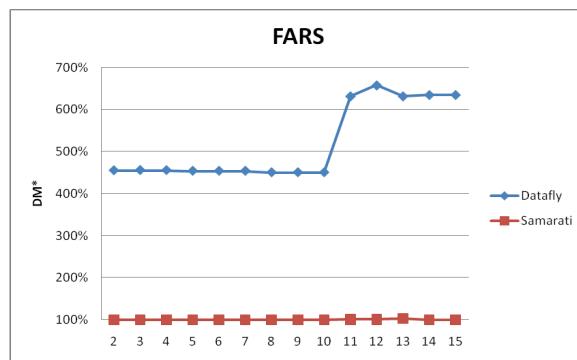


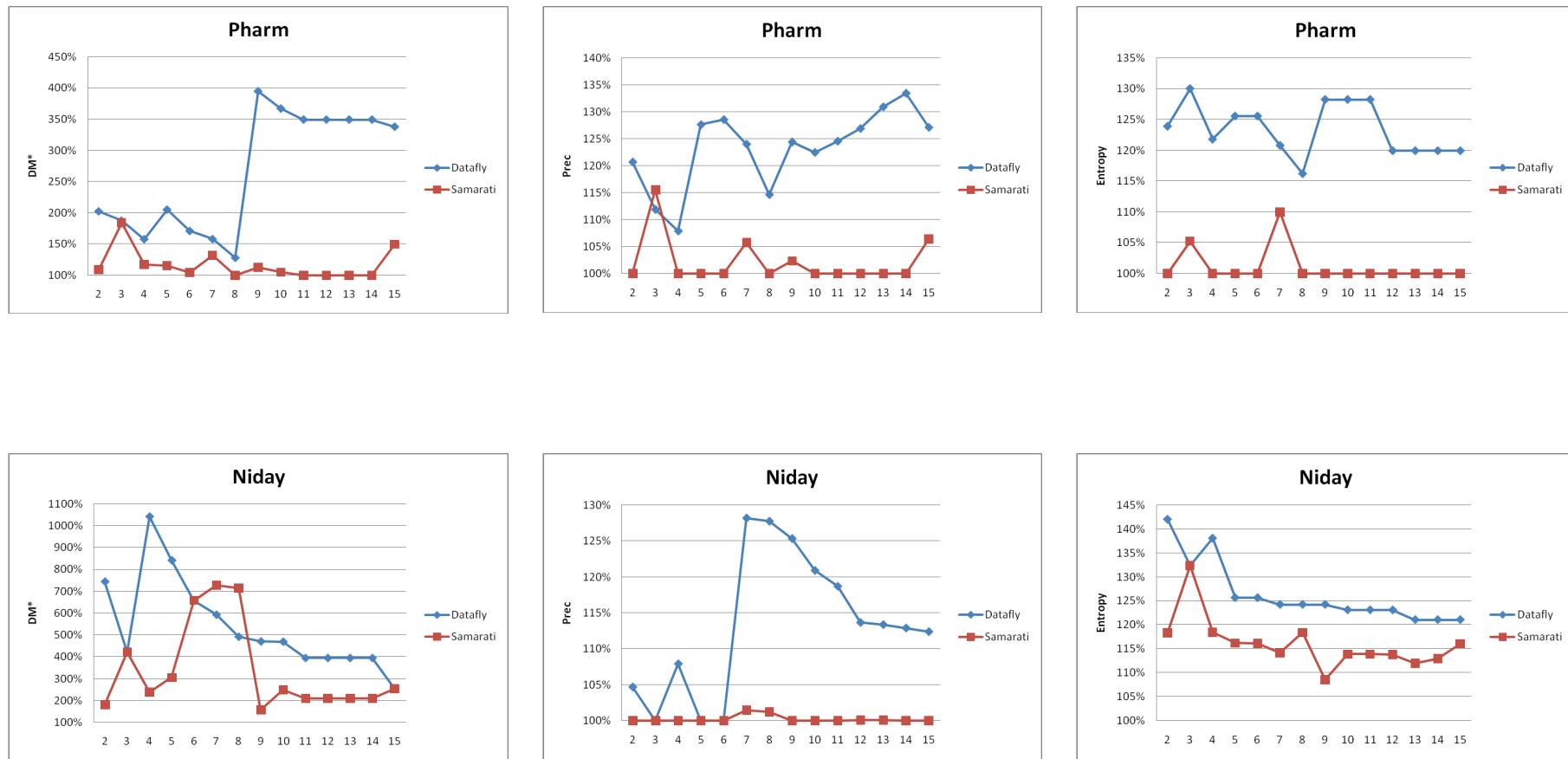




**Figure 1:** Results summary for the six data sets. The x-axis shows the value of  $k$  as it was varied from 2 to 15. The graphs show the information loss results for a 1% suppression limit for the three information loss metrics:  $DM^*$ ,  $Prec$ , and non-uniform entropy.

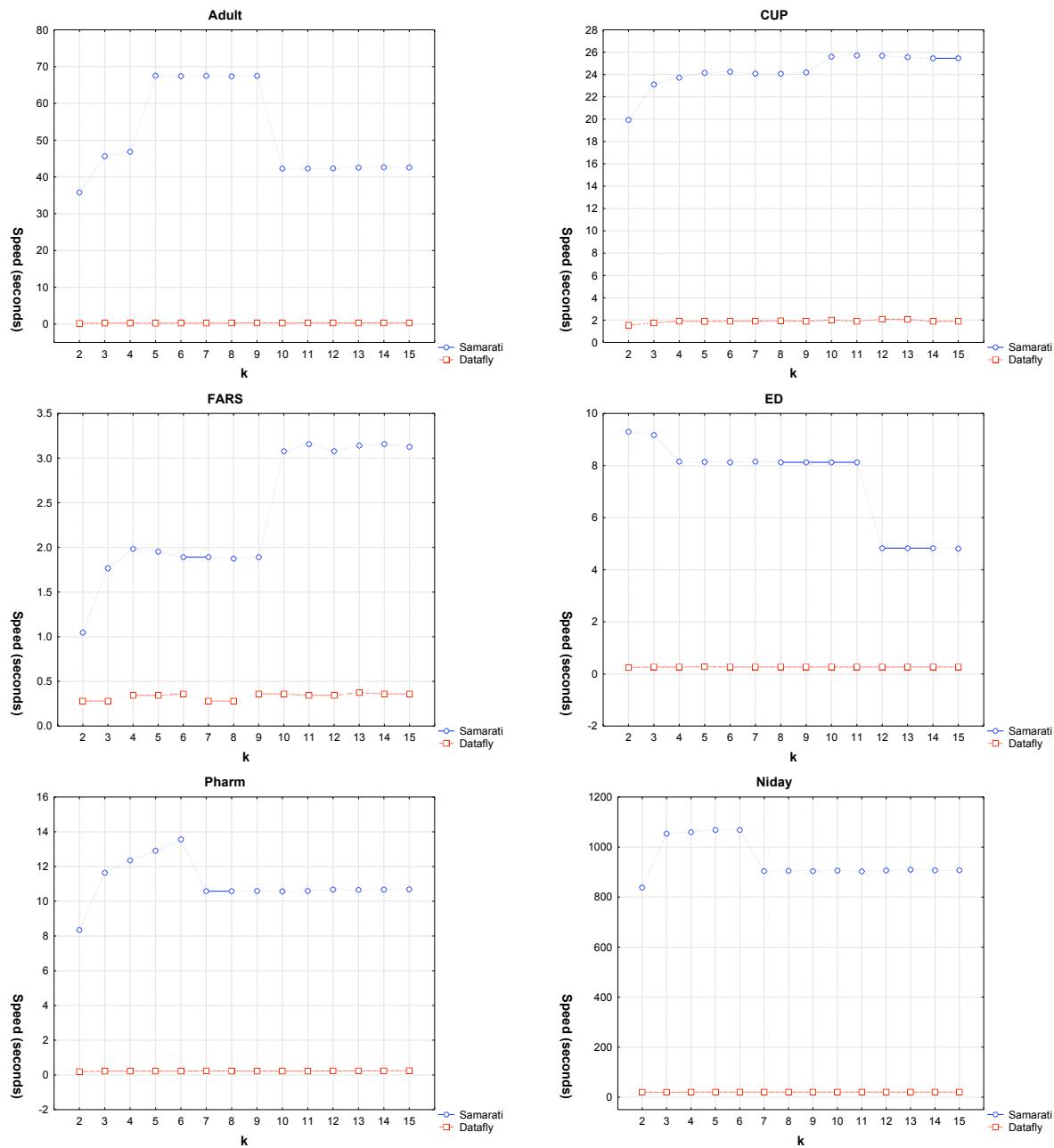
**Information Loss Results for 10% Suppression Limit**





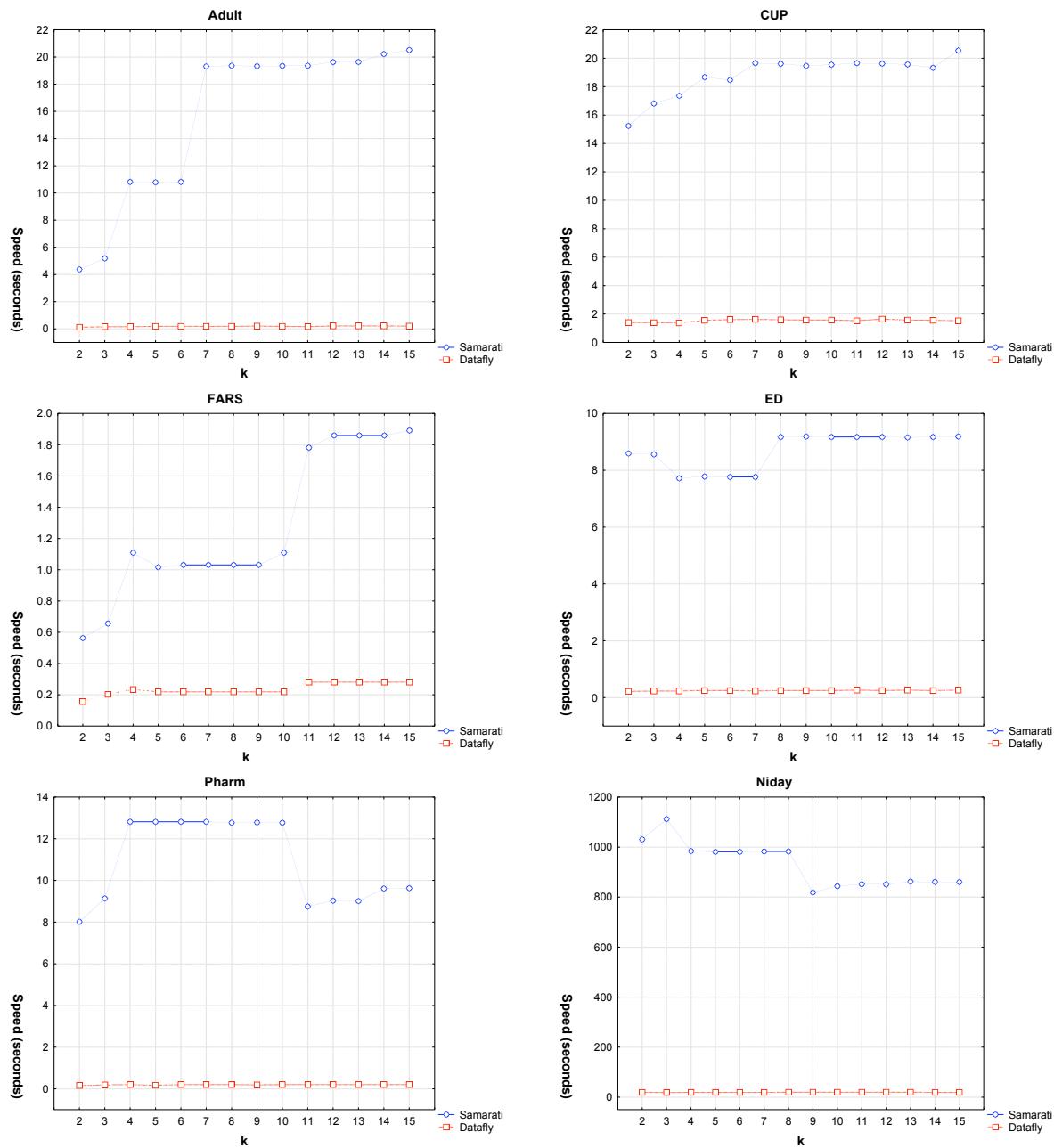
**Figure 2:** Results summary for the six data sets. The x-axis shows the value of  $k$  as it was varied from 2 to 15. The graphs show the information loss results for a 10% suppression limit for the three information loss metrics:  $DM^*$ ,  $Prec$ , and non-uniform entropy.

### Comparison of Datafly vs. Samarati Speed at 1% Suppression Limit



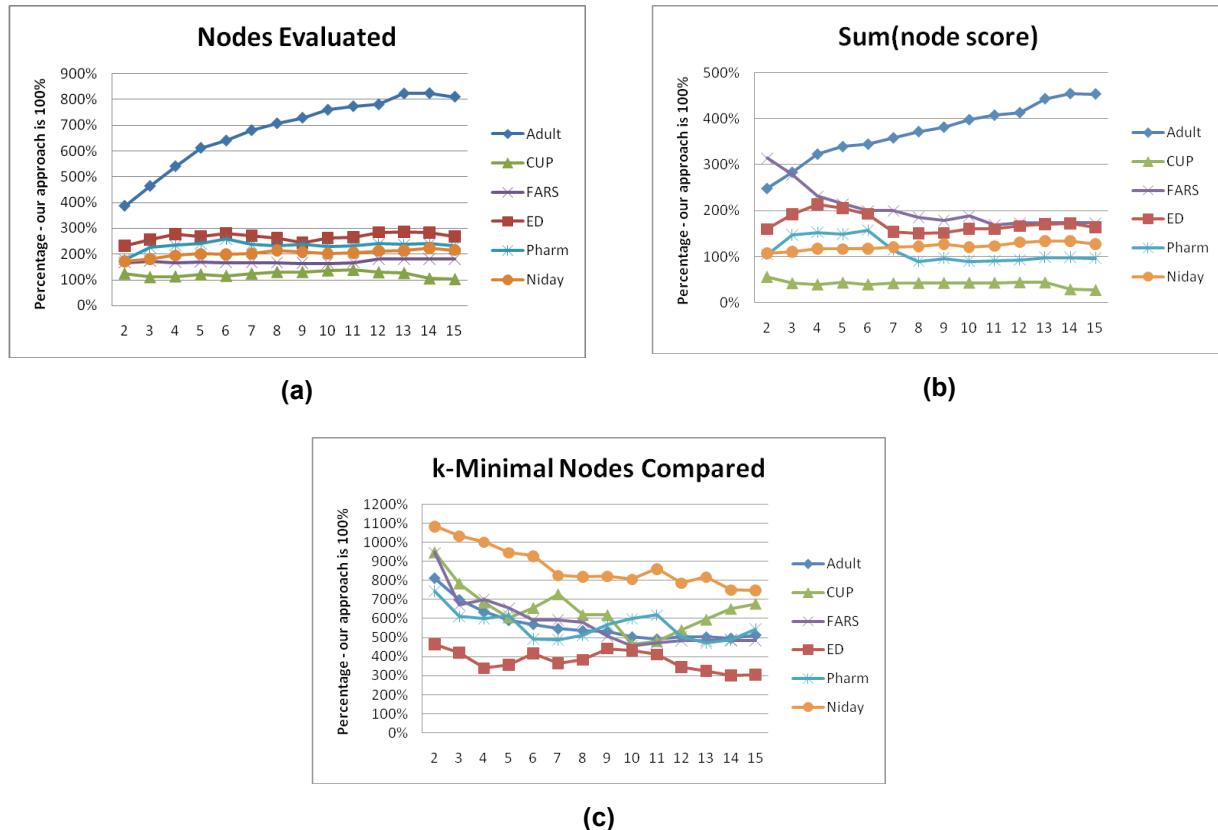
**Figure 3:** The speed to obtain a solution for Datafly and Samarati on the six data sets across values of  $k$  2 to 15 and a maximum suppression value of 1%. The timing is in seconds and includes all pre- and post-processing that may be needed to run each algorithm. These results were obtained on a dual core 3.2 GHz Pentium D processor with 3GB of RAM.

### Comparison of Datafly vs. Samarati Speed at 10% Suppression Limit



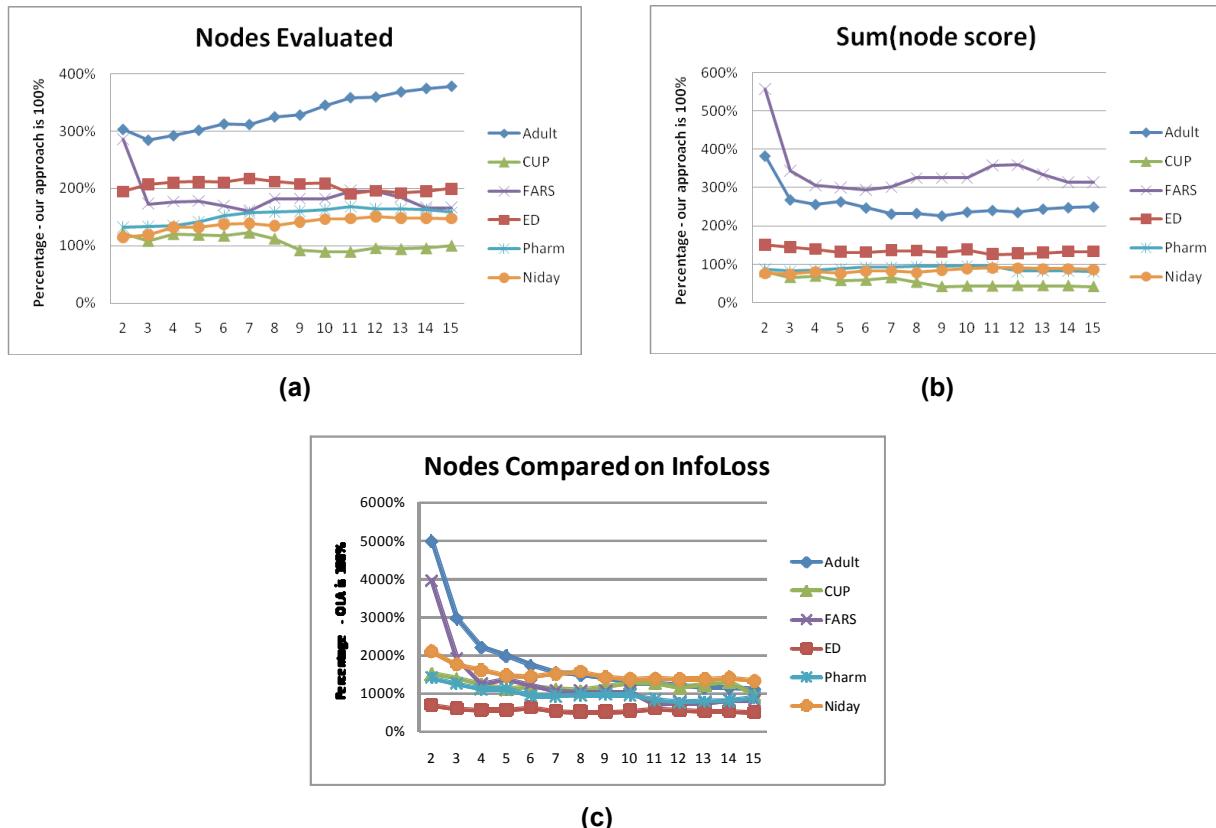
**Figure 4:** The speed to obtain a solution for Datafly and Samarati on the six data sets across values of  $k$  2 to 15 and a maximum suppression value of 10%. The timing is in seconds and includes all pre- and post- processing that may be needed to run each algorithm. These results were obtained on a dual core 3.2 GHz Pentium D processor with 3GB of RAM.

### Comparison of Computations at 1% Suppression Limit



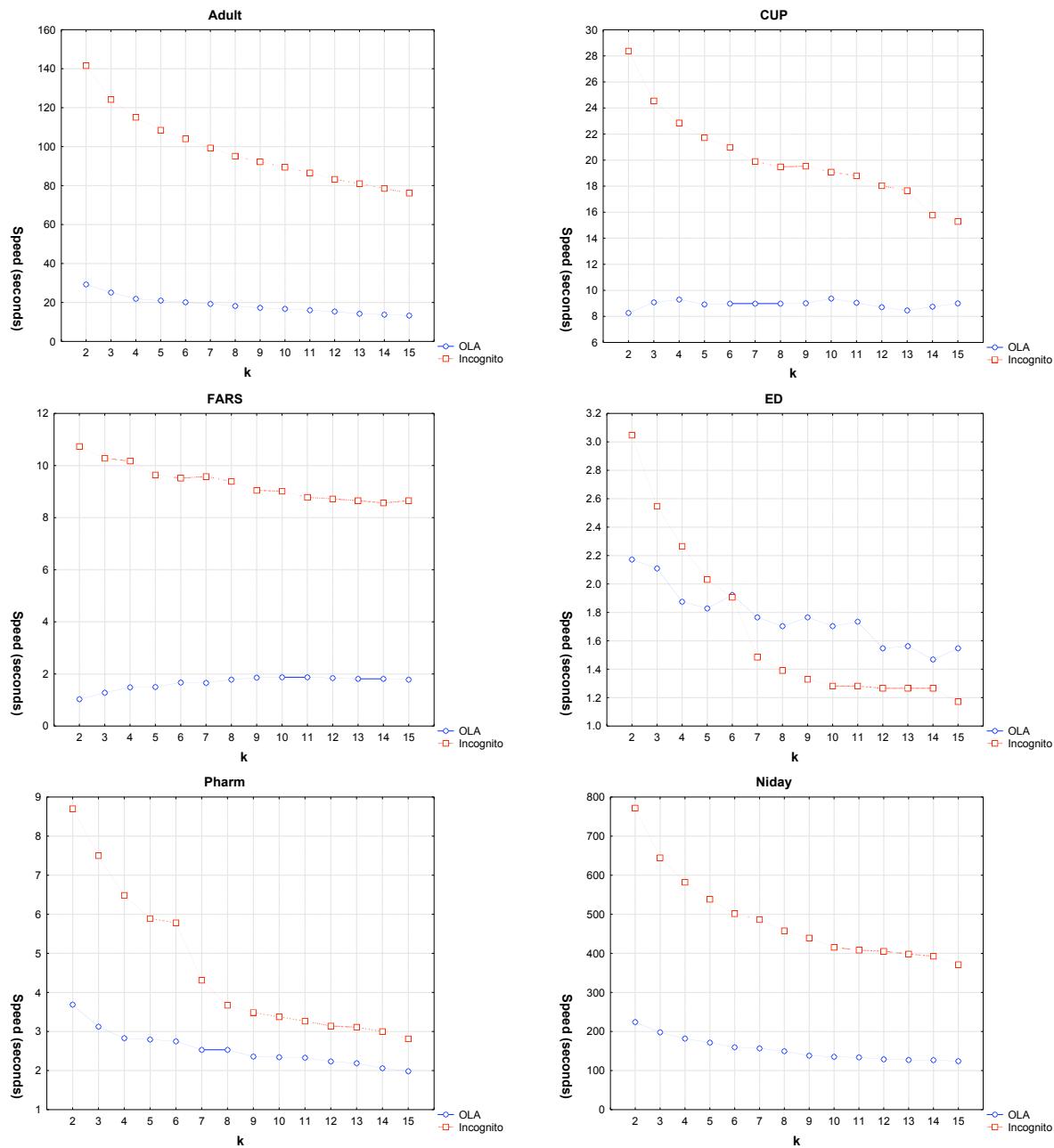
**Figure 5:** The performance metrics comparing our algorithm to Incognito. The results are for the 1% suppression limit. Our algorithm is the 100% value on the y-axis, and if Incognito performs more computations then its value is above 100%, and if it performs less computation then its value is below 100%. The panels show: (a) the total number of nodes for which we need to compute if they are k-anonymous, (b) the node complexity score given by equation (3), and (c) the number of nodes for which information loss needs to be computed.

### Comparison of Computations at 10% Suppression Limit



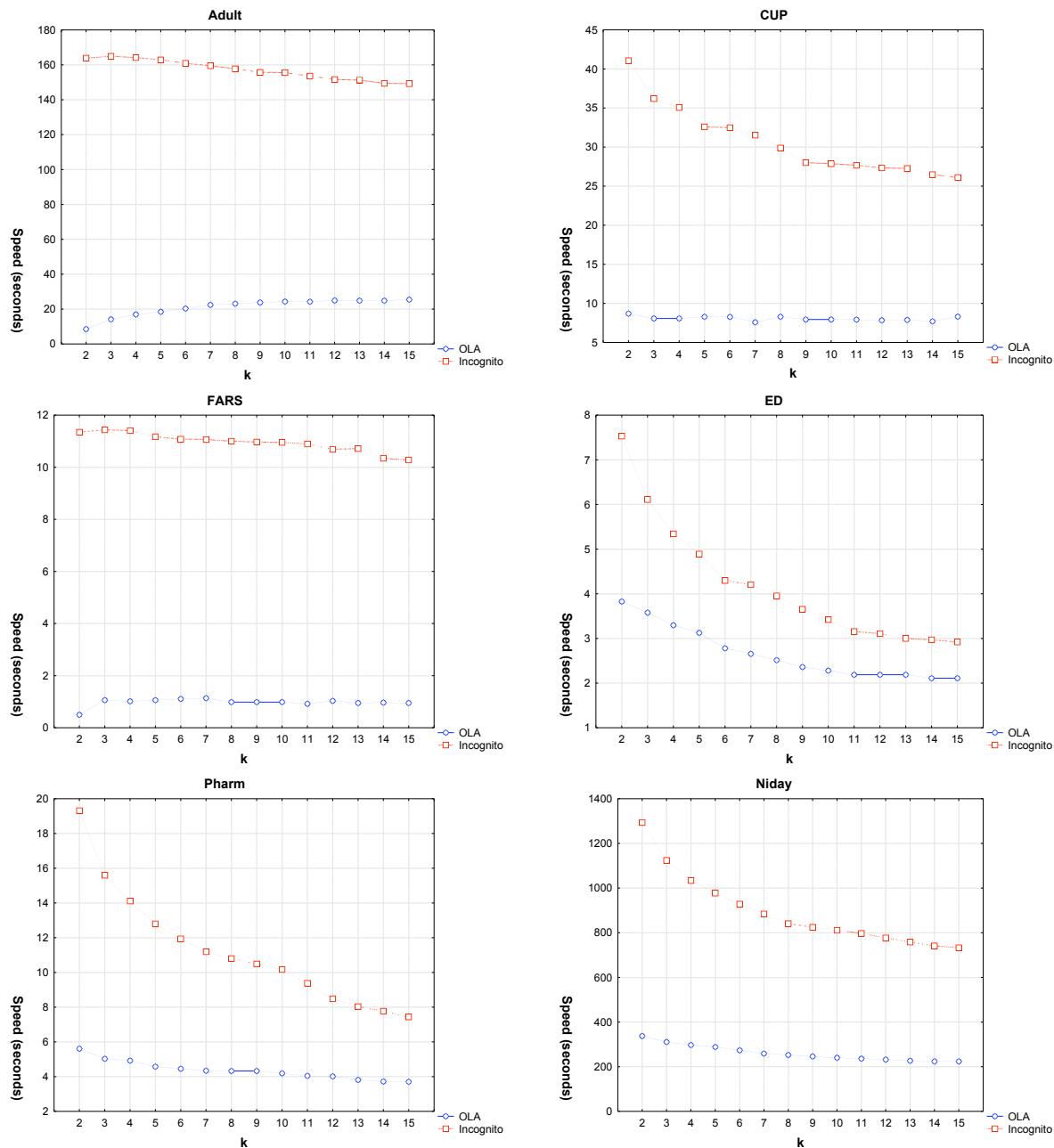
**Figure 6:** The performance metrics comparing our algorithm to Incognito. The results are for the 10% suppression limit. Our algorithm is the 100% value on the y-axis, and if Incognito performs more computations then its value is above 100%, and if it performs less computation then its value is below 100%. The panels show: (a) the total number of nodes for which we need to compute if they are k-anonymous, (b) the node complexity score given by equation (3), and (c) the number of nodes for which information loss needs to be computed.

### Comparison of OLA vs. Incognito Speed at 1% Suppression Limit



**Figure 7:** The speed to obtain a globally optimal solution using OLA and Incognito on the six data sets across values of  $k$  2 to 15 and a maximum suppression value of 1%. The timing is in seconds and includes all pre- and post-processing that may be needed to run each algorithm. These results were obtained on a dual core 3.2 GHz Pentium D processor with 3GB of RAM.

### Comparison of OLA vs. Incognito Speed at 10% Suppression Limit



**Figure 8:** The speed to obtain a globally optimal solution using OLA and Incognito on the six data sets across values of  $k$  2 to 15 and a maximum suppression value of 10%. The timing is in seconds and includes all pre- and post-processing that may be needed to run each algorithm. These results were obtained on a dual core 3.2 GHz Pentium D processor with 3GB of RAM.