

# Population visualization tool for Lambeth Council

Final Project Report

Author: Jan Aldous Torres

Supervisor: Dr Simon Miles

Student ID: 1323454

April 15, 2017

### **Abstract**

(Will be filled after the completion of the report). The abstract is a very brief summary of the report's contents. It should be about half-a-page long. Somebody unfamiliar with your project should have a good idea of what your work is about by reading the abstract alone.

### **Originality Avowal**

I verify that I am the sole author of this report, except where explicitly stated to the contrary.  
I grant the right to King's College London to make paper and electronic copies of the submitted work for purposes of marking, plagiarism detection and archival, and to upload a copy of the work to Turnitin or another trusted plagiarism detection service. I confirm this report does not exceed 25,000 words.

Jan Aldous Torres

April 15, 2017

### **Acknowledgements**

I would like to thank my supervisor Dr. Simon Miles for his support in this project. It is much appreciated.

I thank Lambeth Council for their support in the requirements and evaluation stages of the project. Taking time off their busy schedules was much appreciated.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	Report Structure . . . . .	5
<b>2</b>	<b>Background</b>	<b>6</b>
2.1	Data analysis in government . . . . .	6
2.2	Customer segmentation . . . . .	6
2.3	Related applications . . . . .	7
2.4	Data visualization and exploration . . . . .	8
<b>3</b>	<b>Requirements</b>	<b>10</b>
3.1	Objective . . . . .	10
3.2	Description of data to be used . . . . .	10
3.3	Functional Requirements . . . . .	11
3.4	Creating groups . . . . .	12
3.5	Questions to visualize . . . . .	13
3.6	Non-functional requirements . . . . .	14
<b>4</b>	<b>Specification &amp; Design</b>	<b>15</b>
4.1	System architecture . . . . .	15
4.2	User interface . . . . .	16
<b>5</b>	<b>Report Body</b>	<b>17</b>
5.1	Section Heading . . . . .	17
<b>6</b>	<b>Implementation</b>	<b>18</b>
6.1	Section Heading . . . . .	18
<b>7</b>	<b>Professional and Ethical Issues</b>	<b>19</b>
7.1	Section Heading . . . . .	19
<b>8</b>	<b>Results/Evaluation</b>	<b>20</b>
8.1	Software Testing . . . . .	20
8.2	Section Heading . . . . .	20
<b>9</b>	<b>Conclusion and Future Work</b>	<b>21</b>
	Bibliography . . . . .	22

<b>A</b>	<b>Extra Information</b>	<b>23</b>
A.1	Tables, proofs, graphs, test cases, ... . . . . .	23
<b>B</b>	<b>User Guide</b>	<b>24</b>
B.1	Instructions . . . . .	24
<b>C</b>	<b>Source Code</b>	<b>25</b>
C.1	Instructions . . . . .	25

# List of Figures

4.1	Pie chart visualization . . . . .	15
-----	-----------------------------------	----

# Chapter 1

## Introduction

[?] [5] [6] [7] [?] [9] [1] [2] [8] [4] [3]

Local government councils are challenged with the recent cuts in local government funding and a diverse population with different needs. Therefore, there is a need to efficiently allocate resources. The usual approach is to tackle issues raised by individual social services. However, there has been a shift in local government, which is promoted by the national government, to utilize the concept of customer segmentation. This will let them understand the needs of the population better thus letting them utilize and allocate resources in the most effective way based on those needs [4]. This approach would require the council to identify and therefore focus resources on vulnerable population groups. For example, instead of focusing on issues on education, they would focus on the children in care who are most susceptible to lower educational, social and employment outcomes. Some social services deal with the same population groups so this approach would require an integration of data from the different services. This gives an overview of all data from services provided by the councils. This is unlike the former approach which requires individual social services to conduct data analysis using their own data.[TJA1]

Lambeth Council in London wants to implement the new approach. However, current methods to analyze and visualize data take a long time before the analyses and visualizations would reach the commissioners, the people in the council who formulate strategies to tackle issues.

There are off-the-shelf software packages available which supports users in customer segmentation from clustering to the visualization of their data. In government, using data analysis in policy decision making is becoming more popular(new york guy video). The Local Government



Agency has created a guideline to support councils in tailor-made clustering and visualization [4]. Moreover, there have been projects by other local councils to visualizations of data on prominent population groups.

Lambeth Council would like to have a bespoke population segmentation tool that would be customized to their own needs and available data. This project aims to implement a prototype of such a tool as a Django web application.

The tool will focus on the situation in which population groups are in through visualizations. The situations they face will be defined by the data. Visualizations of this data will also help to compare characteristics of the group to the wider population and subgroups. This comparison may lead to observations on which problems a group or subgroup is facing. This may lead to the conclusion that the council may want to increase resources spent on a group or service. The tool basically aid the commissioner to apply customer segmentation on their population and identify which problems segments of the population are facing thus help make decision on the allocation of their resources.

Clustering will provide a more complete perspective into the different situations groups face instead of viewing a group as one entity and defines variables which distinguish the group from the rest of the population using clustering algorithms. Data integration from open government datasets will create a more in depth view of groups. And possible relations between geographical locations’s characteristics ie facilities (GP), situation (crime) in a ward and whether they are on the London living wage. Since the user may not necessarily know what they are looking for when they read the data, the visualizations will be in the form that will let the user explore more aspects of the data.

## 1.1 Report Structure

The remainder of this paper is organized as follows. Chapter 2 introduces related work to the tool and data exploration, and Chapter 3 discusses the requirements of the tool. Chapter 4 outlines the specification and design of the tool. Chapter 5

## Chapter 2

# Background

This chapter includes work done on formulating better ways of with data exploration. It will later describe in more detail the work done on the access and visualization of open government data to non-technical people who require access to this type of data. The chapter ends with related application to the prototype of this project.

- incorporates all relevant issues

- main issues are analysed and developed

### 2.1 Data analysis in government

Give overview of using data in policy decision making

There has been work to make government data more accessible to people which do not have the technical abilities to merge and manipulate data such as journalists, data journalists, etc. [10]. This is especially true in open government data platforms where datasets take non-standard forms and whose data may be encoded in a way that is not immediately understandable. There have been efforts to evaluate government platforms and to enumerate features that would effectively let users make sense of the data [8].

### 2.2 Customer segmentation

Market segmentation or customer segmentation is concept used in business regarding the division of a homogenous population into segments which have similar attributes, wants, needs or demands. Its objective is to design a marketing mix that precisely matches the expectations of customers in the targeted segment. Few companies are big enough to supply the needs

of an entire market; most must breakdown the total demand into segments and choose those that the company is best equipped to handle. (businessdictionary.com). The four basic market segmentation-strategies are based on: behavioral, demographic, psychographic, and geographical differences (businessdictionary.com). Much of the need to have this is due to the lack of resources needed for the whole market. For many of the local government units, it is the constraint of the budget that would require an efficient allocation of resources by choosing specific groups. The national government and LGA promote the application of this concept usually used in the private sector into local government.

Clustering algorithms could be used to divide the population.

According to the LGA [1], the collation of demographic data and the use of clustering algorithms such as k-means could be used to divide the population into segments. Choosing which

## 2.3 Related applications

There have been applications developed similar to the tool this project aims to implement. These include more general purpose customer segmentation tools which allows the user to analyze and then segment data using algorithms, to the visualization and use of that data. There have also been local governments which are following the concept of customer segmentation to also spend their resources more effectively.

Customer segmentation tool like Mosaic by Experian (mosaic) and Acorn by CACI (acorn) are tools both to allow population segmentation through customer profiling based on set of demographic and other indicators.

Customer classification is a concept promoted by the government and the Local Government Association (LGA) has created a guidance document for any council that wishes to implement such a tool (LGA) (smart cities).

Kent and Medway has tool which segments its population by social class and aims to highlight the key features which make each Group distinctive, to help you visualise the segmentation data and understand the essence of each Group (Kent Medway). It lets the user compare the values between all groups. This tool is bespoke to Kent Medway's data and the tool cannot be reused for other council's data.

There have been similarities between tools and recommendations to differentiate and describe each group. (Kent) (LGA) suggested to include maps showing the concentration of that group in each ward and in addition to a textual description, they include pictures to describe

each group. (Barnet) has a report on the customer segmentation of its population, and includes the percentage and textual description of its population. (Kent) includes graphs to visualize data but also word maps of key characteristics. (LGA) has suggested that the use of spider diagrams which detail the variables compared to the town and district average. Similarly (Acorn) (Kent), they graph pieces of data with an index of the group compared to the overall population. This visualizes both the average value and the groups relation to that average displays whether they are higher or lower than that average however each graph is created for each variable.

LGA has suggested data integration to enable other uses of data. (smart cities) defined use of explicit, "records of who has or is using a service" and implicit "knowledge of staff on customers using a service" customer data in an effort to know their customers in great detail.

## 2.4 Data visualization and exploration

Numerous applications exist to allow users to interact with data. Data exploration involves different types of interaction and functionality and there has been work to create a framework for data exploration [9].

There has been work done on enumerating the problems with data exploration through visualizations [10].

There has also been work done on the visualization of clusters through graphs, which show how closely related the clusters are to each other [7].

Show thinking of why some features of related applications are used here Use cases/features of Kent and Medway: - See group data, more specifically see pictures and phrases which describe the group, distribution between wards, data about services used, benefits used and demographics of the group as an index value (0-200) and a map of the group's population density in the council's area Features of neighbourhood.statistics.gov.uk - The user selects the data to be displayed on the map and chart. Map and geographical unit breakdown chart interacts with each other where hovering over one will highlight the area on both map and chart. The map is divided into geographical units (i.e. ONS's lower output areas) which are colored according to the band which the area is in. The chart is according to the user's setting. - CSV of Lambeth's 2016 Residential Survey: a survey of a sample of Lambeth's population consisting of 1024 people about their quality of life, what they thought of Lambeth's services and about the respondents themselves. The data is in the form of

CSV text files where each row contains a person's entries as categorical data in the form of code (see below Survey Code Translation). Each single answer question has its own column and entries are in code as an integer between 1-100. For the questions which have multiple answers, each answer is regarded as a sub question in the form of `Q5A` meaning question `Q5`, choice code `A` (makes it look like a string but its supposed to be an integer) (see below Survey Code Translation). Therefore, a sub question has its own column and entries are in code as a `Q1A` or `Q0A`. There are also other fields which have been added to the original survey results such as group, subgroup, quintile, which is a result of previous data analysis.

CSV of Lambeth's 2016 Residential Survey: a survey of a sample of Lambeth's population consisting of 1024 people about their quality of life, what they thought of Lambeth's services and about the respondents themselves. The data is in the form of CSV text files where each row contains a person's entries as categorical data in the form of code (see below Survey Code Translation). Each single answer question has its own column and entries are in code as an integer between 1-100. For the questions which have multiple answers, each answer is regarded as a sub question in the form of `Q5A` meaning question `Q5`, choice code `A` (makes it look like a string but its supposed to be an integer) (see below Survey Code Translation). Therefore, a sub question has its own column and entries are in code as a `Q1A` or `Q0A`. There are also other fields which have been added to the original survey results such as group, subgroup, quintile, which is a result of previous data analysis.

Lambeth's 2016 Residential Survey Code Translation: a Microsoft Word document of the original survey with the original questions and code used in the data associated with the question. Under each question is the choice code and English meaning of the choice. The choice code is either a number for single answer questions (e.g. 1. Male, 2. Female) or letters for multiple answer questions (e.g. A. Access to nature, B. Activities for teenagers, etc.). The system may be able to save access (through URLs to CSV and JSON files) or actual files of the datasets (CSV or JSON files themselves). The URLs must be kept in persistent storage. If the input is in the form of a file, the file must be stored. The system may be able to save access (through URLs to CSV and JSON files) or actual files of the datasets (CSV or JSON files themselves). The URLs must be kept in persistent storage. If the input is in the form of a file, the file must be stored."

## Chapter 3

# Requirements

The functional and non-functional requirements was created based on communication, through interviews and email, with a member of the Policy and Communications Team, and requirements which I have created based on initial requirements by our contact. Our contact described the context in which the tool may be used, the overall objective of the software, the users' technical abilities and some functional requirements, namely the method to create groups (see section 3.4). Communication with our contact was not sustained after initial requirements were collected due to the council's large workload. The rest of the requirements was created in accordance with the overall goal, initial requirements and according to my own analysis of the data.

### 3.1 Objective

The purpose of this application is to support those in local government councils, in charge of creating policies, to explore data on residents demographics and satisfaction with council services. User-defined groups will allow the user to focus on specific portions of the population. Aided by visualizations through graphs and maps, it may lead the user to identify which segments of the population requires the which resources.

### 3.2 Description of data to be used

The following local data provided by Lambeth should be used:

- **CSV of Lambeth's 2016 Residential Survey:** a survey of a sample of Lambeth's population consisting of 1024 people about their quality of life, what they thought of

Lambeth's services and about the respondents themselves. The data is in the form of a CSV text file. Residents' answers are categorical data in the form of code specified by the Survey Code Translation (see point below, Survey Code Translation). Single answer questions have its own column and its choice code is a positive integer. Multiple answer questions, each answer is regarded as a sub question in the form of 'Q5A' meaning question 5, choice code 'A' (see below Survey Code Translation). Therefore, a sub question has its own column and entries are in code as a 1 or 0 meaning yes or no respectively. There are also other fields which have been added to the original survey results such as group, subgroup, quintile, which is a result of previous data analysis.

- **Microsoft Word document of Lambeth's 2016 Residential Survey Code Translation:** a Microsoft Word document of the original survey with the original questions and code used in the data associated with the question. Under each question is the choice code and English meaning of the choice. The choice code is either a number for single answer questions (e.g. 1. Male, 2. Female) or letters for multiple answer questions (e.g. A. Access to nature, B. Activities for teenagers, etc.).
- **GeoJSON of Lambeth's Ward's Boundary Specifications:** downloaded from Lambeth's website, it is a GeoJSON, simple geographical features encoded as a JSON, of the geographical boundaries of each ward.

### 3.3 Functional Requirements

Functional requirements with the words `textbfshould` or `textbfmust` is a required feature. Requirements with the word `textbfmay`, is a desirable feature which may or may not be implemented should there not be enough development time.

1. Data storage
  - (a) 1. The system must save the actual files of the datasets (CSV, text files or JSON) and not access the file through a URL of another website.
2. Grouping requirements
  - (a) The user should be able to create groups based on the parameters to the 5 factors (see section 3.4). A group's information should be kept in a database.
  - (b) The user should be able to view data about a group through graphs.

- i. Single answer questions or questions which requires only one answer should be visualized as stacked bar charts or pie charts.
  - ii. Multi-code questions or questions which requires more than one answer should be visualized as bar charts.
  - iii. Clicking data on a graph should display which wards have answered the selected question and choice on the map (see functional requirement 4).
  - iv. The questions under section Questions to be visualized should be visualized following requirements 3a and 3b.
- (c) Like a population density map, which shows higher density concentrations of people as a darker color and lower density concentrations as a lighter color on particular areas of the map, the map should show the selected data (see functional requirement 3c). The map should be divided into the council's wards; the boundaries of each ward should be obvious. The color of the ward should depend on the number of residents who have the selected question and answer (see functional requirement 3c).
  - (d) The user should be able to compare values of a data variable between all groups such that groups should be compared to the average value of the variable in the whole population (i.e. compare the percentage of disabled people who are male to the percentage of the whole population who are male).
  - (e) 1. The main data used in the visualizations should be the CSV of Lambeth's Residential Survey and its Code Translation.

### 3. Clustering

- (a) The system should segment a group using a clustering algorithm.
- (b) The system should display any statistical information (e.g. the mean of the answers of a question) on the differences between the clusters.
- (c) The user should be able to compare information between the group's data and each of the clusters' data.

## 3.4 Creating groups

A group should be based on residents' answers to 5 questions in the residential survey. At least 1 question is required to form a group, therefore the rest of the 4 questions could have any



answer. The group will be composed of residents who match the required answer(s) to the required questions.

There are 5 questions in which a user can create a group:

1. Whether they have a disability or long term illness (Q43)
2. What type of benefits do they acquire (Q45A-Q45K)
3. Educational/employment activity (Q46)
4. Whether they are on the London Living Wage (Q47)
5. Housing tenure (i.e. council tenant, private owner, etc.) (Q35)

\*Question number in the CSV residential survey is in brackets.

### 3.5 Questions to visualize

Once a group has been created, the user must be able to see the answers to the following survey questions:

1. What matters most to them most (Q5) (Multiple, top 3)
2. How was their last contact with the Council was made (Q26A- Q26G) (Multiple)
3. How they use the website (Q29) (Multiple)
4. What services they have used (Q39) (Multiple)
5. How they access the internet (Q50) (Multiple)
6. How well the changes have benefited them (Q11) (Single)
7. Whether they feel they belong to their neighborhood (Q13\_R1) (Single)
8. Whether they value the friendships in their neighborhood (Q13\_R2) (Single)
9. Whether they could approach a neighbor for advice (Q13\_R3) (Single)
10. Whether neighbors help out each other (Q13\_R4) (Single)
11. Whether they would be willing to work with others to improve their neighborhood (Q13\_R5) (Single)
12. Whether they would join community events in their area (Q13\_R6) (Single)

13. Whether they regularly stop and talk to people in their neighborhood (Q13\_R7) (Single)
14. Whether they would speak highly of their neighborhood when asked (Q13\_R8) (Single)
15. Gender (QGEN) (Single)
16. Age (QAGE) (Single)
17. Ethnicity (QETH) (Single)

\*Question number in the residential survey is in brackets followed by whether they are a single answer question (Single) or multiple answer question (Multiple).

The answers to the questions above should be visualized using graphs for each group and their clusters based on functional requirement 3.

### 3.6 Non-functional requirements

**Usability:** The users for the program are the policy makers of the council, those responsible for formulating strategies to allocate resources for a council based on the presentation of data analysis given to them. They do not necessarily have technical skills, in terms of being able to operate applications, or advanced statistical analysis skills. Since the users are not technical savvy, the ease of use is imperative to the design of the user interface. The visualizations should also be easily interpreted.

In terms of maintainability, this version of the system is only a prototype to show the potentials of such a system, therefore there will not be a need for maintainability of the code. In terms of security, the data being used is anonymous therefore there will be no need for security measures for the data.

## Chapter 4

# Specification & Design

This section describes a high-level design of the system to be implemented. As mentioned in the Requirements section, this will also be the first version of the design documentation. It will be added on to in the future.

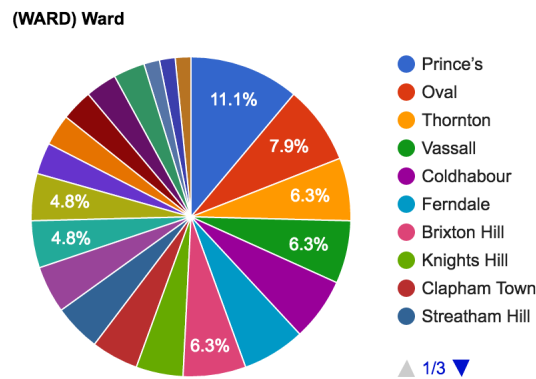


Figure 4.1: Pie chart visualization

### 4.1 System architecture

Due to the iterative nature of implementation, the architecture chosen will make it easy to add more components and even more data sources. The system will follow the three-layer architecture. The data layer will consist of the survey data and cluster criteria, parameters (based on the 5 clustering factors) for clusters set by the user. The logic layer will consist of creating clusters and extracting clusters from survey data. The presentation layer includes the

form to create clusters and the navigation between clusters and visualization of cluster data.

## 4.2 User interface

The system's main goal is to present data and let the user interact with the user in an effective way. Thus the design of the user interface is integral to the success of the system.

The cluster creation UI will enable the user to name and enter the parameters for a cluster. The choices for the parameters will be constrained to the fields in the survey and the field's answers.

To enable the easily navigate between each cluster, there will be a navigation pane in the form of tabs for each cluster. Once a tab is selected, the graphs of data listed in section (data visualization) will be displayed. To enable comparison between data on each cluster, inspired by (kent and midway), each graph will be shown as indexed horizontal bar charts (figure). This will visualize the comparison of the values of the cluster to average value for that field in the whole population.

# Chapter 5

## Report Body

The central part of the report usually consists of three or four chapters detailing the technical work undertaken during the project. **The structure of these chapters is highly project dependent.** They can reflect the chronological development of the project, e.g. design, implementation, experimentation, optimisation, evaluation, etc (although this is not always the best approach). However you choose to structure this part of the report, you should make it clear how you arrived at your chosen approach in preference to other alternatives. In terms of the software that you produce, you should describe and justify the design of your programs at some high level, e.g. using OMT, Z, VDL, etc., and you should document any interesting problems with, or features of, your implementation. Integration and testing are also important to discuss in some cases. You may include fragments of your source code in the main body of the report to illustrate points; the full source code is included in an appendix to your written report.

### 5.1 Section Heading

#### 5.1.1 Subsection Heading

## Chapter 6

# Implementation

### 6.1 Section Heading

## Chapter 7

# Professional and Ethical Issues

Either in a separate section or throughout the report demonstrate that you are aware of the **Code of Conduct & Code of Good Practice** issued by the British Computer Society and have applied their principles, where appropriate, as you carried out your project.

### 7.1 Section Heading

## Chapter 8

# Results/Evaluation

### 8.1 Software Testing

### 8.2 Section Heading



## Chapter 9

# Conclusion and Future Work

The project's conclusions should list the key things that have been learnt as a consequence of engaging in your project work. For example, "The use of overloading in C++ provides a very elegant mechanism for transparent parallelisation of sequential programs", or "The overheads of linear-time n-body algorithms makes them computationally less efficient than  $O(n \log n)$  algorithms for systems with less than 100000 particles". Avoid tedious personal reflections like "I learned a lot about C++ programming...", or "Simulating colliding galaxies can be real fun...". It is common to finish the report by listing ways in which the project can be taken further. This might, for example, be a plan for turning a piece of software or hardware into a marketable product, or a set of ideas for possibly turning your project into an MPhil or PhD.

# References

- [1] 3. customer segmentation and profiling.
- [2] Customer segmentation.
- [3] Customer segmentation.
- [4] Developing a customer classification tool: Guidance document for local authorities.
- [5] Kent & medway interactive guide.
- [6] Brigitte Boden, Roman Haag, and Thomas Seidl. Detecting and exploring clusters in attributed graphs: a plugin for the gephi platform. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, CIKM '13, pages 2505–2508. ACM.
- [7] Drew Conway. Data science through the lens of social science. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, pages 1520–1520. ACM.
- [8] Anna Gogolou, Marialena Kyriakidi, and Yannis Ioannidis. Data exploration: A roll call of all user-data interaction functionality. In *Proceedings of the Third International Workshop on Exploratory Search in Databases and the Web*, ExploreDB '16, pages 31–33. ACM.
- [9] Alvaro Graves and James Hendler. Visualization tools for open government data. In *Proceedings of the 14th Annual International Conference on Digital Government Research*, dg.o '13, pages 136–145. ACM.

# Appendix A

## Extra Information

### A.1 Tables, proofs, graphs, test cases, ...

The appendices contain information that is peripheral to the main body of the report. Information typically included in the Appendix are things like tables, proofs, graphs, test cases or any other material that would break up the theme of the text if it appeared in the body of the report. It is necessary to include your source code listings in an appendix that is separate from the body of your written report (see the information on Program Listings below).

## Appendix B

# User Guide

### B.1 Instructions

You must provide an adequate user guide for your software. The guide should provide easily understood instructions on how to use your software. A particularly useful approach is to treat the user guide as a walk-through of a typical session, or set of sessions, which collectively display all of the features of your package. Technical details of how the package works are rarely required. Keep the guide concise and simple. The extensive use of diagrams, illustrating the package in action, can often be particularly helpful. The user guide is sometimes included as a chapter in the main body of the report, but is often better included in an appendix to the main report.

# Appendix C

## Source Code

### C.1 Instructions

Complete source code listings must be submitted as an appendix to the report. The project source codes are usually spread out over several files/units. You should try to help the reader to navigate through your source code by providing a “table of contents” (titles of these files/units and one line descriptions). The first page of the program listings folder must contain the following statement certifying the work as your own: “I verify that I am the sole author of the programs contained in this folder, except where explicitly stated to the contrary”. Your (typed) signature and the date should follow this statement.

All work on programs must stop once the code is submitted to KEATS. You are required to keep safely several copies of this version of the program and you must use one of these copies in the project examination. Your examiners may ask to see the last-modified dates of your program files, and may ask you to demonstrate that the program files you use in the project examination are identical to the program files you have uploaded to KEATS. Any attempt to demonstrate code that is not included in your submitted source listings is an attempt to cheat; any such attempt will be reported to the KCL Misconduct Committee.

**You may find it easier to firstly generate a PDF of your source code using a text editor and then merge it to the end of your report. There are many free tools available that allow you to merge PDF files.**