

# Clustering and Fitting (40%)

Start Assignment

- Due 25 Apr by 12:00
- Points 40
- Submitting a file upload
- File types pdf
- Available 22 Mar at 14:00 - 30 Apr at 12:00

<b>Weighting %:</b>	40	<b>Submission deadline (for students):</b>	25/4/24 at 12pm (Midday)
<b>Authorship:</b>	Individual	<b>Target date for returning marked coursework:</b>	17/5/24
<b>Tutor setting the work:</b>	Dr. William Cooper	<b>Number of hours you are expected to work on this assignment:</b>	20

## **This Assignment assesses the following module Learning Outcomes (from Definitive Module Document):**

1. Be able to maintain and develop code using the git version control system.
2. Be able to apply different techniques for cleaning data and preparing it for analysis.
3. Be able to design and implement algorithms for clustering, classification and regression problems.
4. Be able to communicate their findings to others, including a critical assessment of performance.
5. Demonstrate knowledge and understanding of the concepts of version control for code development.
6. Demonstrate knowledge and understanding of key data manipulation techniques for data preparation.
7. Understand how to approach a range of different data science problems to obtain an efficient solution.

## **Assignment Tasks:**

You will create a well-written report performing clustering and fitting within a dataset. You can download any dataset from Kaggle/Worldbank/etc. Be sure to include your name, student number and a link to your GitHub repository in the report. There will be at least four plots: a histogram/bar chart/pie chart; a line/scatter graph; a confusion matrix/heatmap/corner/box/violin plot; an elbow/silhouette plot. The code will contain evidence of the creation of any displayed graphs (one graph per function) and the creation of any shown clustering/fitting technique. The minimum expected techniques will be that of k-means clustering and line fitting.

This will build on the statistics and trends assignment into a full report as would be produced by a professional data scientist. However, do **not** use the same report/dataset as previously (this will be checked), as self-plagiarism is still academic misconduct.

## **Submission Requirements:**

A three page PDF report, including a functional link to your GitHub repository containing your python code (either notebooks or plain python). Check that your repository link is both clickable and links to a **public** repository.

**Marks awarded for:**

See rubric.

**Type of Feedback to be given for this assignment:**

Written feedback within the rubric.

**Additional information:**

- Regulations governing assessment offences including Plagiarism and Collusion are available from [https://www.herts.ac.uk/\\_data/assets/pdf\\_file/0007/237625/AS14-Apx3-Academic-Misconduct.pdf](https://www.herts.ac.uk/_data/assets/pdf_file/0007/237625/AS14-Apx3-Academic-Misconduct.pdf) (https://www.herts.ac.uk/\_data/assets/pdf\_file/0007/237625/AS14-Apx3-Academic-Misconduct.pdf) (UPR AS14) .
- Guidance on avoiding plagiarism can be found here: <https://herts.instructure.com/courses/61421> (https://herts.instructure.com/courses/61421) (see the Referencing section)
- For postgraduate modules:
  - a score of 50% or above represents a pass mark.
  - late submission of any item of coursework for each day or part thereof (or for hard copy submission only, working day or part thereof) for up to five days after the published deadline, coursework relating to modules at Level 7 submitted late (including deferred coursework, but with the exception of referred coursework), will have the numeric grade reduced by 10 grade points until or unless the numeric grade reaches or is 50. Where the numeric grade awarded for the assessment is less than 50, no lateness penalty will be applied.

**Assignment 2: Clustering and Fitting**

Criteria	Ratings			Pts
<b>Relational Graph Quality</b> The quality of the relational graph, e.g. line/scatter graph.	<b>2 Pts</b> <b>Full marks</b> The graph will convey an xy relation. The axes labels will be fully readable without effort and the relation(s) will be clear.	<b>1 Pts</b> <b>Fair quality</b> The graph will convey an xy relation. The axes labels may be too small to read comfortably. There may be an overcrowding of the figure.	<b>0 Pts</b> <b>No marks</b> Missing graph or missing axes labels.	2 pts
<b>Categorical Graph Quality</b> The quality of the categorical graph, e.g. bar chart/histogram/pie chart.	<b>2 Pts</b> <b>Full marks</b> The graph will compare multiple categories. The axes labels will be fully readable without effort and the appearance will be clear.	<b>1 Pts</b> <b>Fair quality</b> The graph will compare multiple categories. The axes labels may be too small to read comfortably. There may be an overcrowding of the figure.	<b>0 Pts</b> <b>No marks</b> Missing graph or missing axes labels.	2 pts
<b>Statistical Graph Quality</b>	<b>2 Pts</b> <b>Full marks</b>	<b>1 Pts</b> <b>Fair quality</b>	<b>0 Pts</b> <b>No</b>	

The quality of the statistical graph, e.g. heatmap/confusion matrix/corner plot/violin/box plot.	The graph will communicate a statistical relation. The axes labels will be fully readable without effort and the appearance will be clear.			The graph will communicate a statistical relation. The axes labels may be too small to read comfortably. There may be an overcrowding of the figure.			<b>marks</b> Missing graph or missing axes labels.	2 pts
Quality of Analysis How accurate and meaningful the data analysis is.	<b>5 Pts Full marks</b> The explanation is clear and coherent. Statistics are used to support statements. There is a connecting storyline.	<b>4 Pts Very high marks</b> The explanation is clear and coherent. Statistics are used to support statements. There may be some connecting storyline.	<b>3 Pts High marks</b> The explanation is mostly clear and coherent. There may be some statistics supporting some statements. There may be some storyline.	<b>2 Pts Fair marks</b> The explanation is mostly coherent. There may be a majority of statements without statistical support. The report is more descriptive.	<b>1 Pts Poor quality</b> The report is almost entirely descriptive without meaningful statistics.	<b>0 Pts No marks</b> No description of any merit.	5 pts	
Spelling and Grammar The quality of the overall use of English.	<b>1 Pts Good</b> The spelling and grammar is acceptable enough to communicate complex ideas.		<b>0.5 Pts Acceptable</b> The spelling and grammar use is acceptable enough to communicate basic ideas.			<b>0 Pts No marks</b> Very poor English, making idea communication challenging.		1 pts
Relational Graph Function The function in the code that creates the relational graph.	<b>1 Pts Good</b> Function with docstring which only creates one plot.		<b>0.5 Pts Acceptable</b> Function without docstring or function produces multiple plots.			<b>0 Pts No marks</b> No/not useable GitHub link or no function.		1 pts
Categorical Graph Function The function in the code that creates the categorical graph.	<b>1 Pts Good</b> Function with docstring which only creates one plot.		<b>0.5 Pts Acceptable</b> Function without docstring or function produces multiple plots.			<b>0 Pts No marks</b> No/not useable GitHub link or no function.		1 pts
Statistical Graph Function The function in the code that creates the statistical graph.	<b>1 Pts Good</b> Function with docstring which only creates one plot.		<b>0.5 Pts Acceptable</b> Function without docstring or function produces multiple plots.			<b>0 Pts No marks</b> No/not useable GitHub link or no function.		1 pts
Statistical Depth	<b>3 Pts</b>		<b>2 Pts</b>		<b>1 Pts</b>	<b>0 Pts</b>		

The depth of the statistics used in the code.	<b>Full marks</b> All major moments shown (mean/median, standard deviation, skewness, kurtosis). Correlation matrix and basic 'describe' used.			<b>High marks</b> First two major moments shown (mean/median, standard deviation). Correlation matrix and basic 'describe' used.		<b>Fair marks</b> Correlation matrix and basic 'describe' used.	<b>No marks</b> No/not useable GitHub link or no function or no use of 'describe' and correlation matrix.	3 pts
Code Quality The appearance of the code and adherence to PEP-8.	<b>2 Pts Full marks</b> Code is easy to read and follows the major PEP-8 recommendations: import > functions > variables order; functions separated by exactly two lines (one if in a class) or sole occupier of notebook cell; spaces after commas and around assignment/mathematical operators.		<b>1 Pts Fair marks</b> Code is mostly easy to read and may have a few slips from the major PEP-8 recommendations: import > functions > variables order; functions separated by exactly two lines (one if in a class) or sole occupier of notebook cell; spaces after commas and around assignment/mathematical operators.		<b>0 Pts No marks</b> No/not useable GitHub link or code is difficult to read with many divergences from the major PEP-8 recommendations: import > functions > variables order; functions separated by exactly two lines (one if in a class) or sole occupier of notebook cell; spaces after commas and around assignment/mathematical operators.			2 pts
Clustering Function The function in the code that performs the clustering.	<b>1 Pts Good</b> Function with docstring which does not create a plot.			<b>0.5 Pts Acceptable</b> Function without docstring or function also creates plots.		<b>0 Pts No marks</b> No/not useable GitHub link or no function.		1 pts
Fitting Function The function in the code that performs the fitting.	<b>1 Pts Good</b> Function with docstring which does not create a plot.			<b>0.5 Pts Acceptable</b> Function without docstring or function also creates plots.		<b>0 Pts No marks</b> No/not useable GitHub link or no function.		1 pts
Clustering Quality How well the clustering has been performed.	<b>6 Pts Full marks</b> The clusters will appear well grouped. The data will have been normalised and back scaled to present. Clear use of silhouette score/elbow method to select	<b>5 Pts Very high marks</b> The clusters will appear well grouped. The data will have been normalised. Clear use of silhouette score/elbow method to select cluster amount.	<b>4 Pts High marks</b> The clusters will appear well grouped. The data may have been normalised. Use of silhouette score/elbow method to select cluster amount. The graph	<b>3 Pts Fair marks</b> The clusters will appear well grouped. The data may have been normalised. The graph will have coloured groups and labelled cluster centres in the legend. The data	<b>2 Pts Poor quality</b> The clusters may appear well grouped. The data may have been normalised. The data will be appropriate for clustering.	<b>1 Pts Very poor quality</b> The clusters are not well grouped. The data may be appropriate for clustering.	<b>0 Pts No marks</b> The clusters are not well grouped. The data is not appropriate for clustering, or no graph.	6 pts

	cluster amount. The graph will have coloured groups and labelled cluster centres in the legend. The data will be appropriate for clustering.	The graph will have coloured groups and labelled cluster centres in the legend. The data will be appropriate for clustering.	will have coloured groups and labelled cluster centres in the legend. The data will be appropriate for clustering.	will be appropriate for clustering.					
Fitting Quality How well the fitting has been performed.	<b>5 Pts Full marks</b> The data will be well fitted and suitable for fitting. The plot will include a good confidence interval and errorbars.	<b>4 Pts High marks</b> The data will be well fitted and suitable for fitting. The plot will include a good confidence interval or errorbars.	<b>3 Pts Fair marks</b> The data will be well fitted and suitable for fitting. The plot may include reasonable errorbars.	<b>2 Pts Poor quality</b> The data will be fitted and suitable for fitting.	<b>1 Pts Very poor quality</b> The data will be poorly fitted but suitable for fitting.	<b>0 Pts No marks</b> The data is not suitable for fitting, or no graph.	5 pts		
Clustering Prediction Accuracy of clustering predictions.	<b>4 Pts Full marks</b> Several predicted points will be attached to appropriate groups, and clearly labelled and coloured.	<b>3 Pts High marks</b> Several predicted points will be attached to appropriate groups.	<b>2 Pts Fair marks</b> Predictions will be made for different group memberships.	<b>1 Pts Poor quality</b> An attempt at predictions will have been made for different group memberships.		<b>0 Pts No marks</b> No predictions made or no/not useable GitHub link.	4 pts		
Fitting Prediction Accuracy of fitting predictions.	<b>3 Pts Full marks</b> Several predictions with good, associated uncertainties are shown.		<b>2 Pts Fair marks</b> Several predictions are given.	<b>1 Pts Poor quality</b> An attempt at predictions are made.		<b>0 Pts No marks</b> No predictions made or no/not useable GitHub link.		3 pts	
Submission Guidelines Keeping normal text and margins whilst maintaining the expected page length.	<b>0 Pts Expected</b> The report is at the required length with no overly small text or minimised margins.		<b>-4 Pts Not expected</b> The report is not at the required length, either overrunning or too short (by at least a third of a page). Alternatively, there may be overly small text or minimised margins.					0 pts	

