

SCC403 – Data Mining

Coursework Assignment

1 Introduction

A data scientist must be able to process various data sets and streams, to know various methods and techniques and be able to select and apply the most suitable algorithms for data mining.

In particular, techniques such as data pre-processing, data clustering, and classification.

The objective of this assignment is to conduct data analysis across two tasks. The first task tackles a climate dataset and the second task centres around Deep Neural Networks (DNNs) on a selection of datasets. The assignment includes selection and justification of the specific methods for data pre-processing (normalisation, standardisation, feature selection/extraction, anomaly detection (if any), missing data (if any)), their implementation and analysis of the results as well as a well annotated code. You are expected to critically analyse the results of applying these techniques, and demonstrate a clear understanding of the purposes and processes of data analysis. In addition to your report, please submit your source code, including comments. To achieve top marks, a well justified variety of specific techniques and well commented code are expected. Analysis and understanding of the methods, algorithms and the overall process are the most important elements in addition to the implementation skills such as the code and the presentation.

We expect the use of Python - the most widely used language for machine learning which we also use in the labs, but if you prefer to use a different language we may need to contact you for clarification, if we believe that your code is not running correctly.

2 Task 1: Basel Climate Dataset

You are expected to use the set of climate data provided in the file '*ClimateDataBasel.csv*'. This data is a subset of publicly available (from <https://www.meteoblue.com/>) data about climate in Basel, Switzerland which contains 1763 18-dimensional records of data from the summer and the winter seasons of the period from 2010 to 2019. The meaning of each column of data is listed below:

- Temperature (Min) °C.
- Temperature (Max) °C.
- Temperature (Mean) °C.
- Relative Humidity (Min) %.
- Relative Humidity (Max) %.
- Relative Humidity (Mean) %.
- Sea Level Pressure (Min) hPa.
- Sea Level Pressure (Max) hPa.

- Sea Level Pressure (Mean) *hPa*.
- Precipitation Total *mm*.
- Snowfall Amount *cm*.
- Sunshine Duration *min*.
- Wind Gust (Min) *Km/h*.
- Wind Gust (Max) *Km/h*.
- Wind Gust (Mean) *Km/h*.
- Wind Speed (Min) *Km/h*.
- Wind Speed (Max) *Km/h*.
- Wind Speed (Mean) *Km/h*

2.1 Preprocessing

Before clustering, it's essential to preprocess the Basel Climate Dataset to ensure clean and consistent data. You should research and apply appropriate techniques, which may include:

- Normalization/Standardization
- Outlier Detection
- Feature Selection/Extraction
- Dealing with Missing Data

Justify your preprocessing choices in your report. You may provide techniques outside of this list, especially when aiming to achieve top marks.

2.2 Clustering

The objective is to cluster the climate data set. Choose at least two clustering algorithms and apply them to the climate data set. To achieve top marks one of the methods should be from independent research.

Develop the programme and explain the functionality of the algorithms in as much detail as you can. Compare the results and limitations of each of the algorithms that you have used. Provide qualitative and quantitative evaluation of the results and visualise them.

3 Task 2: Image Processing with Deep Neural Networks (DNNs)

In this task, you will explore the use of pre-trained Deep Neural Networks (DNNs) for image processing. Using these models allows us to transform complex data (e.g., images) into lower-dimensional, information-dense representations known as feature vectors. These vectors form a so called *latent space* or feature space in which similar objects (e.g. images) are mapped near each other while dissimilar ones are separated further apart. Your goal is to use one of these models to extract feature vectors for further analysis, including visualization, clustering, and classification.

You are required to select **two** datasets from the following options:

- OxfordPetsIIIT (<https://www.robots.ox.ac.uk/~vgg/data/pets/>)

- iRoads (<https://www.cs.auckland.ac.nz/~m.rezaei/Publications/iROADS%20Dataset.pdf>)
- Cats vs Dogs (<https://www.kaggle.com/datasets/karakaggle/kaggle-cat-vs-dog-dataset>)
- Food101 (<https://www.kaggle.com/datasets/dansbecker/food-101>)

3.1 Preprocessing/Feature Extraction

Preprocess each selected dataset as necessary (e.g., resizing images, normalization, etc.) to ensure it is compatible with the provided DNN models. The preprocessing step is crucial for obtaining meaningful feature representations.

Choose one, or more, of the DNN models and extract feature vectors from each chosen dataset:

- DinoV2 (<https://github.com/facebookresearch/dinov2>)
- VGG16 (<https://pytorch.org/vision/stable/models/generated/torchvision.models.vgg16.html>)
- DenseNet (<https://pytorch.org/vision/stable/models/densenet.html>)
- ResNet (<https://pytorch.org/vision/main/models/resnet.html>)

Choose a technique (t-SNE, PCA, UMAP, see Lab 4) to reduce the dimensionality of your feature vectors for visualisation and clustering and visualise in 2D.

3.2 Clustering

Perform clustering on the selected data sets within the *latent* feature space using a suitable clustering algorithm (e.g., K-Means, ELM, K-Medoids). Evaluate the quality of the clusters using purity, Davies-Bouldin or other metrics found in literature (see also Lab 5 and Lecture 5)

3.3 Classification

Apply classifier(s) to the feature vectors. For example, add a linear fine-tuneable layer(s) to the DNN used for feature extraction as in Lab 9. Provide a detailed discussion on its performance, including accuracy, confusion matrices, and any other relevant metrics.

When analysing the performance of the classifiers you should use precision/recall, F1 score and classification accuracy. You may also indicate the time required for training the classifier as a measure of computational complexity (note that the time is always conditional on the type of hardware you use - laptop, computer, CPU/GPU, etc.) and is not an absolute measure, but when making comparisons it can be useful.

Hints: You may also compare the performance of the DNN with added linear layer with other classification methods. You are encouraged to explore recent research papers to identify advanced classification techniques that are well-suited for working with high-dimensional feature spaces, for example using a Nearest Neighbour approach in regards to the clustered data in the latent space.

4 Marking Scheme

The marks are allocated as follows:

Report (20%):

- Structure and presentation (6%)
- Language and style (7%)
- Use of literature and references (7%)

Task 1 (40%):

- Data Pre-processing (total 20%)
- Clustering (total 20%)

Task 2 (40%):

- Data Pre-processing (total 15%)
- Clustering (total 10%)
- Classification (total 15%)

Each subsection of Tasks 1 and 2 (e.g., Data Pre-processing, Clustering, and Classification) will be evaluated based on the following criteria, with each criterion carrying equal weight:

- Level of understanding
- Depth of analysis
- Working, well annotated code and results
- Justification of selected methods
- Independent research and use of methods not given in the lectures

At the end of this document there is an Appendix, which explains what a mark means in Lancaster University and includes suggestions for a well-written report.

The length of the report should not exceed 6 pages. You can use double column format, e.g. the so-called IEEE style as described in the Appendix. You may include an Appendix (4 pages maximum) following the main report.

5 Deadlines and general requirements

The lectures and tutorials will provide you with the necessary tools to conduct your analysis. You may also include additional analysis methods that you have researched separately, that may help derive your conclusion (this is not compulsory).

You are expected to critically analyse the results of applying these techniques, and demonstrate a clear understanding of the purpose and processes of data analysis. **The deadline for submission is: Friday 13th December 2024 by 4:00PM.** The cut-off deadline is Monday 16th December 2024 by 4:00PM (with late submission penalty incurred which is 1 letter grade or 10%). Submissions after this deadline cannot be accepted according to the University regulations.

In case your code is unclear to us you may be contacted for interview. If you fail to reply or attend the interview your code could be marked as “not working”.

6 Additional Comments

You must report in an “acknowledgements” section the use of any libraries, readily available online code, and code from online tutorials. Additionally, you are free to discuss your work with colleagues, but you must also report in the “acknowledgments” section if anyone has helped you significantly. Remember that using others’ work without giving the due credit is an act of *plagiarism*, and it is not a good academic practice.

APPENDIX

Example of the style of the report

Title of the Report

Student number

line 1: dept. name of organization

line 2-name of the programme and module

Abstract— Briefly describe the outline of your report. You can download a template (Word or LaTeX from <https://www.ieee.org/conferences/publishing/templates.html>)

I. Introduction

Here you have to provide the background review. of the existing approaches stressing the ones that have been actually used. Critically analyse and compare alternative techniques and methods. Try to go beyond what was given in the lectures using external sources and references.

II. Pre-processing

Here you have to provide a description and description and the results of pre-processing techniques that are relevant and stress those that you actually used in your work. Provide the software code that you used to obtain the results in an Appendix. Do not forget to justify your choice.

III. Clustering

Here you have to provide a description and the results of clustering techniques that are relevant and stress those that you actually used in your work. Provide the software code that you used to obtain the results in an Appendix. A very important part of your report is the critical analysis of the results. Why have you used the stated methods? What are the advantages and limitations of the algorithms that you have used?

IV. Classification

Here you have to provide a description and the results of classification techniques that are relevant and stress those that you actually used in your work. Provide the software code that you used to obtain the results in an Appendix. A very important part of your report is the critical analysis of the results. Why have you used the stated methods? What are the advantages and limitations of the algorithms that you have used?

V. Conclusion

Describe briefly what has been done, with a summary of the main results. Discuss here possible future developments (what you would have done more). What is distinctive about the results you have obtained?

VI. References

The template will number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in [3]—do not use “Ref. [3]” or “reference [3]” except at the beginning of a sentence: “Reference [3] was the first ...”

Number footnotes separately in superscripts. Place the actual footnote at the bottom of the column in which it was cited. Do not put footnotes in the reference list. Use letters for table footnotes.

- [1] J. Han, M. Kamber, Data Mining: Concepts and Techniques, Morgan Kaufmann, 2001
- [2] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. Heidelberg, Germany: Springer Verlag, 2001
- [3] Angelov, P.: Autonomous Learning Systems: From Data Streams to Knowledge in Real Time. John Wiley and Sons (2012).
- [4] Angelov, P.: Outside The Box:An Alternative Data Analytics Framework. *Journal of Automation, Mobile Robotics & Intelligent Systems*. Vol. 8, 29–35.

Appendix

Please include here additional experimental results or additional details.

Presenting someone else’s work as your own in an assignment without proper citation of the source is an act of **plagiarism**. More information about Lancaster University Plagiarism Framework can be found at <https://www.lancaster.ac.uk/academic-standards-and-quality/information-and-resources/policies-and->

What a Mark Means in Lancaster University

70 + (Distinction)

Critical Understanding of Topic

Excellent understanding and exposition of relevant issues; insightful and well informed, clear evidence of independent thought; good awareness of nuances and complexities; appropriate use of theory.

Structure of Research

Substantial evidence of well implemented independent research and / or Substantial evidence of well selected evidence to support argument.

Use of Literature

Excellent use of literature to support argument /points.

Conclusion

Excellent; clear implications for theory and/or practice.

Language

Excellent; a delight to read.

Structure and Presentation

Arguments clearly structured and logically developed; sensible weighting of parts; meaningful diagrams; properly formatted references.

65 – 69% (Very Good Pass)

Critical Understanding of Topic

Clear awareness and exposition of relevant issues; some awareness of nuances and complexities but tendency to simplify matters; based on appropriate choice and use of theory.

Structure of Research

Some evidence of independent research reasonably well implemented and / or some evidence of identification of suitable evidence to support argument.

Use of Literature

Good use of literature to support arguments.

Conclusion

Very good; draws together main points; some implications for theory and/or practice

Language

Carefully written; negligible errors.

Structure and Presentation

Arguments clearly structured and logically developed; good weighting of parts; meaningful diagrams; properly formatted references.

60 – 65% (Good Pass)

Critical Understanding of Topic

Shows awareness of issues and theories; attempts at analysis but tendency to lapse into description

Structure of Research

Some evidence of independent research reasonably well implemented and / or some evidence of identification of suitable evidence to support argument.

Use of Literature

Use of standard literature to support arguments.

Conclusion

Reasonable conclusion that summarises essay; a few implications for theory and/or practice.

Language

A few errors; generally satisfactory.

Structure and Presentation

Arguments reasonably clear but undeveloped; some meaningless diagrams or poor structure.

50 – 59% (Pass)

Critical Understanding of Topic

Work shows understanding of topic but at superficial level; no more than expected from attendance at lectures; some irrelevant material; too descriptive.

Structure of Research

Insufficient evidence of independent research and / or very limited evidence used to support argument.

Use of Literature

Use of secondary literature to support arguments.

Conclusion

Conclusion does not do justice to body of essay; too short; no implications.

Language

Some errors; grammar and syntax need attention.

Structure and Presentation

Arguments not very clear; poor organisation of material; poor use of diagrams; poor referencing.

45 – 49% (Marginal Fail)

Critical Understanding of Topic

Establishes a few relevant points but superficial and confused; much irrelevant material; very little or no understanding of the issues raised by the topic or topic misunderstood; content largely irrelevant; no choice or use of theory; essay almost wholly descriptive; no grasp of analysis with many errors and/or omissions.

Structure of Research

No evidence of independent research and / or No attempt to identify suitable evidence to support argument.

Use of Literature

Relies on a superficial repeat of class notes.

Conclusion

No recognisable conclusion.

Language

Frequent errors; needs urgent attention.

Structure and Presentation

Arguments often confused and undeveloped; no logical structure; very poor organisation of material; many meaningless diagrams; negligible referencing.

0 – 44% (Clear Fail)

Critical Understanding of Topic

Establishes a few relevant points but superficial and confused; much irrelevant material; very little or no understanding of the issues raised by the topic or topic misunderstood; content largely irrelevant; no choice or use of theory; essay almost wholly descriptive; no grasp of analysis with many errors and/or omissions.

Structure of Research

No evidence of independent research and / or No attempt to identify suitable evidence to support argument.

Use of Literature

No significant reference to literature.

Conclusion

No recognisable conclusion.

Language

Frequent errors; needs urgent attention.

Structure and Presentation

Arguments often confused and undeveloped; no logical structure; very poor organisation of material; many meaningless diagrams; negligible referencing.