



# MACQUARIE University

*Faculty of Science & Engineering*

## COMP8325 Applications of Artificial Intelligence for Cyber Security

### Assignment (Part II) Description

#### LEARNING OUTCOME

This assignment (10% weighting) deals with application of machine learning for cybersecurity application. On successful completion, you will be able to:

- Engage with material learned in COMP8325;
- Explain the basic concepts and the limitations of Artificial Intelligence;
- Detect intrusion in networks and systems by applying tools and techniques revealing abnormal patterns in datasets; and
- Analyse the trends of applications of Artificial Intelligence in cyber security.

#### TASK 1: Touch Biometrics

Continuous authentication methods for users who are using a smartphone. A touch dataset (in **data** folder, available at Google Drive [https://drive.google.com/drive/folders/1esR5KU-sw7rqE0i5Sj079\\_0kn5vb8b88?usp=sharing](https://drive.google.com/drive/folders/1esR5KU-sw7rqE0i5Sj079_0kn5vb8b88?usp=sharing)) consists of 30 features (in **features** folder) extracted from raw touch input. Further details about the touch dataset can be found at <http://www.mariofrank.net/touchalytics/index.html>. Please answer the following questions:

- (1) Implement two more features in addition to the 30 found in the database. Do they have positive information gain? That is, are the features useful?
- (2) Report correlation of these feature to the rest of the implemented features.
- (3) Train your model on a binary classifier of your choice (“true user” or “false user” classification problem) using the following 4 scenarios in which you use a feature selection method to choose top 10 features. Describe this process. Use 10-fold cross validation to compute precision and recall in the following scenarios. Try to maximize F1 score when optimizing your classifier. Report F1 and any methods you used to optimize your classifier.
  - (a) 10 top features;
  - (b) 10 top features & your features;

- (c) 30 computed features;
- (d) 30 computed features & your features and qualitatively describe which family of features are most discriminating in your classifier.

## TASK 2: Merits of Entropy in Attack Detection/Diagnostics

Consider a server-log dataset (in `server-log.txt`) hosted at Google Drive [https://drive.google.com/drive/folders/1esR5KU-sw7rqE0i5Sj079\\_0kn5vb8b88?usp=sharing](https://drive.google.com/drive/folders/1esR5KU-sw7rqE0i5Sj079_0kn5vb8b88?usp=sharing). Two attacks happened on a day, both somewhere around 8am and noon. Please answer the following questions:

- Identify the exact date and time<sup>1</sup>. What approach did the attackers use?
- There has been significant literature<sup>23</sup> discussing how entropy can be used to detect these attacks. To do it effectively, approximation schemes are usually used. You do not have to implement these approximation techniques, but do present an analysis of whether entropy is useful and which combinations you tried, e.g. `src ip`, `dest ip`, `src-port`, `dst-port`, etc. Do any reveal anomalies when the two attacks happen?

## SUBMISSION

Please submit your code and analysis (with critical description of the results) in iPython notebook printed as PDF.

## EXPECTATION AND TIMELINE

- Students should submit a single word or pdf file.
- The assignment is due **Monday, 07 June 2021, 9am**.
- Late submissions will incur the following penalties:
  - 10% penalty for 1 to 24 hours late;
  - 20% penalty for 24 to 48 hours late; and
  - 100% penalty for over 48 hours late (iLearn assignment submission automatically closes).
- If you have a legitimate reason for submitting late, discuss this with the convenor well in advance of the assignment due date.

## MARKING RUBRIC

Marks will be available in iLearn by one weeks after the submission due date.

- All the required data analysis tasks have been reasonably accomplished.
- The organisation, presentation and readability of the report
- Appropriate justification of which you have chosen and what you have done in the data analysis process, as well as critical thinking and understanding on the related aspects of the machine learning methods
- The quality of source code, especially the ease of using the code to perform prediction/testing on the reserved data.
- The prediction results in the testing stage (based on the reserved data sets).

---

<sup>1</sup>The columns are labelled.

<sup>2</sup>Lall, et al 2013. *Data Streaming Algorithms for Estimating Entropy of Network Traffic*,

<sup>3</sup>Clifford, Cosma, 2013. *A simple sketching algorithm for entropy estimation over streaming data*