

Intelligent Systems, 2022/23, written exam, 20 January 2023

All questions count equally. Literature, electronic and communication devices are not allowed. It is allowed to use one sheets of A4 format paper with notes. You can write your answers in English or Slovene. Duration: 90 minutes. A possibility to see grading of your exams will take place on Wednesday, 25 January 2023, at 14:00 in the office of Prof Robnik Šikonja (2nd floor, room 2.06).

1. A data scientist has collected a classification dataset consisting of numerical features to solve a given problem in industrial setting. She intends to test several classification algorithms.
Describe for each of the following four algorithms if and how it would benefit from a preprocessing step where numerical features are discretized into a fixed number of discrete bins:
 - i) neural network, i.e. multilayer perceptron,
 - ii) decision rules,
 - iii) support vector machines,
 - iv) random forest.

v) Argue which algorithm would benefit most.
2. A researcher is trying to predict the likelihood of a customer churning from a telecommunications company using historical data. They have trained a logistic regression and decision tree classifiers and achieved an accuracy of 99% on the training set for both of them. However, when tested on a new, unseen dataset, the models' accuracy drops to 89%.
 - a) What are possible problem with the models?
 - b) Discuss possible benefits of regularization techniques such as ridge or lasso for logistic regression to mitigate lower performance in this scenario.
 - c) Discuss possible benefits of or pruning for decision tree to mitigate lower performance in this scenario.
3. A bank is using a machine learning model to detect fraudulent transactions. The dataset they are using is heavily imbalanced, with only 0.2% of the transactions being labelled as fraudulent.
 - a) What are some of the problems that can arise when training a model on such a dataset?
 - b) Describe two techniques that can be used to address this issue and improve the performance of the model in detecting fraudulent transactions.
 - c) Propose a suitable evaluation metric.
4. A team of engineers is working on a sentiment analysis task, and they have decided to use a large generative language model such as GPT-3 to classify the sentiment of the text.
 - a) What are possible advantages of using a large generative pre-trained language model over a BERT language model in sentiment analysis?
 - b) How can the team fine-tune the model on their specific dataset and ensure its robustness?
 - c) Prepare an example of input, i.e. prompt for the model.
 - d) How can the team evaluate the performance of the model?
5. A student is working on a binary classification task and has trained a gradient boosting model to predict if a customer will churn or not. The student wants to understand how the model is making its predictions and how certain features are impacting the predictions.
 - a) How can the student use an explanation technique such as SHAP to understand feature importance and predictions of the model?
 - b) What are the challenges and limitations of SHAP for model interpretation?
 - c) How can the student evaluate the reliability of the SHAP values and ensure that they are providing a fair and accurate representation of the model's predictions?