

Internship project description

A pilot study for federated learning on oncological data

Jan Ramon

11/2023

1 Motivation and context

Over the last decades, there has been an increasing interest in exploiting data. On the other hand, recently there has also been an increasing awareness of the risks of collecting sensitive data centrally, given the frequency of data leaks, hacking or abuse. The Horizon Europe projects TRUMPET and FLUTE will work towards a platform for secure privacy-preserving federated machine learning where the sensitive data remains with the data owners, and machine learning is performed collaboratively by these data owners by participating in collaborative algorithms which (through the use of differential privacy and/or encryption) generate the desired statistical models but prevents sensitive data from being revealed.

These projects also feature use cases in medicine, in particular the TRUMPET project will study lung cancer clustering and eligibility prediction for radiotherapy for head and neck cancer patients, while the FLUTE project will study prediction and diagnosis of prostate cancer.

Before tackling these use cases with federated learning, this internship will conduct a first, short exploratory study to understand how these medical machine learning problems could be solved using machine learning in a simpler setting with central data.

2 Objectives

The goal of this internship project is to find an adequate machine learning approach to at least one of the TRUMPET use cases.

In particular, the objectives are

- to briefly review the literature on the specific machine learning techniques we will use

- to make a description of the format and structure of the available data, the clinical objectives and the relevant background knowledge
- to learn predictive models with the several machine learning algorithms, to investigate what are the best strategies in terms of feature engineering, model selection, hyperparameters and other relevant choices.
- to prepare the next steps, in particular suggest strategies to achieve statistical privacy and to obtain efficient federated algorithms.

3 Plan

Here is a tentative work plan:

- Understanding the data and the clinical problem. Interacting with medical experts and data experts and systematize the knowledge relevant for the project (4 weeks)
- Machine learning literature study (1 week)
- Preparing the data, engineering features and other preprocessing (2 weeks)
- using a selection of algorithms to learn models (3 weeks)
- Fine tuning models, optimizing hyperparameters, combining strategies (3 weeks)
- Suggesting strategies for statistical privacy (2 weeks)
- Suggesting strategies for federated learning (2 week)
- Completion of the internship report (2 weeks)

The timing (here estimated as 20 weeks) can be adapted according to the personal preferences of the student or the requirements of his school or type of project.

4 Environment

The project will be conducted in the INRIA MAGNET team. The student will collaborate and interact with various other collaborators in the team. It is possible the student will be asked to present progress to the team or to partners in ongoing projects. An application for ZRR access will need to be made to the FSD for students of level master-2.

5 Requirements

We expect the student has a good understanding of basic computer science concepts, e.g., data structures and algorithms, and of statistics / machine learning. Knowledge of Python is required.