# Implementing Real-Time Patient Level Predictions Using PLP Models

Scott Brown, Jon Duke, James Fairbanks, Christine Herlihy, Kausar Mukadam, Jason Poovey, Melissa Rost
Georgia Tech Research Institute, Atlanta GA

**Georgia Tech Research Institute**

**ABSTRACT:** *We built a general framework over the OHDSI patient-level prediction package to publish an API for a predictive analytics model for C. Diff, a deadly bacterial infection. This model makes risk predictions for individual patients based within a precision medicine decision support tool with the ability to toggle off and on the inclusion of potential drugs or treatments. This tool can aid medical practitioners when prescribing drugs with potential complications. With this ability, consideration needs to be given to the quality of the model itself as that effects the quality of the predictions and the inability of the model to predict for new treatments not available in the training data.*

## Background

We used the OHDSI Patient Level Prediction (PLP) package to **extract patient-level data** from the MIMIC dataset comprising approximately 40,000 critical care patients. We were specifically interested in **building a predictive model** for Clostridium difficile, commonly known as C. Diff. This life-threatening infection in the colon is caused by taking antibiotics that kill the good gut bacteria so that bad bacteria can grow unchecked. In addition to the problems for the patient, costs attributable to C. Diff in the US are estimated at $6.3 billion annually and nearly 2.4 million days of inpatient stay annually so understanding the causes and factors involved more thoroughly has the potential to save lives and reduce a costly, resource-intensive problem.

The general PLP pipeline involves **creating a cohort of patients** in the database, **extracting features of interest** present in individual patients, and analyzing the resulting dataset by **building a predictive model**. The OHDSI PLP package to perform these tasks is written in R; however, the third step is implemented by having R call specific models in scikit-learn, a state of the art machine learning library written in Python. In order to facilitate access to the entire scikit-learn library, we use scikit-learn directly for the third step of analyzing the resulting data.

## Integration with the Patient-Level Prediction Package

**Index and outcome cohorts are first created** using either dynamically constructed and executed SQL queries or OHDSI ATLAS. Next, we use the features defined in the PLP package to **extract patient-level information**. In order to support both **bulk queries for training** and **individual patient queries for predicting**, the features were extracted using QueryGarden SQL queries generated for both purposes. The resulting data from the cohorts, population, outcomes, and covariates is used in the analytics step.

Contact: melissa.rost@gtri.gatech.edu







## Predicting Outcome Risk for a Single Patient Using a Trained Model

In order to predict C. Diff risk, we **build a predictive model** using scikit-learn. A web **user interface** (UI) allows researchers to choose any combination of features (i.e., drugs, conditions, procedures, and demographics) and model type (i.e., random forest or logistic regression) for training of this model and prediction.

A web application built on QueryGarden supports feature extraction over a REST API for this model. This extends the PLP pipeline of training a model and evaluating its accuracy by **enabling a real-time system to use the trained model to predict the risk of C. Diff for specific new patients.**

This capability is deployed in both the Advanced Clinical Decision Support (ACDS) system to provide precision medicine and the Patient Risk Viewer to summarize hospital-wide risk. For potential drugs to prescribe, the system **predicts an individualized expected patient outcome** to aid a physician in clinical decision-making.

## Conclusions

Utilizing the OHDSI PLP package, we were able to extend the predictive model built from a cohort of patients in order to **provide a real-time prediction for new patients with the same selection criteria for the cohort**. Such a tool allows health practitioners the ability to gain knowledge of how an individual patient could be expected to react to a new drug or treatment before having to go through prescribing the drug and waiting for it to run its course.

This kind of advancement has the potential to revolutionize how medicine is prescribed; however, the medical community needs to be aware of potential flaws. The predictions are **only as good as the model** used to create them and creating the model is the real 'science' of the data science. Additionally, a predictive model only works with data that it was trained on, so **predicting for new covariates that weren't previously seen yields invalid results**. Precision medicine going forward will require collaboration between data scientists and medical professionals to understand the potential benefits and flaws.