

# Predicting Feasibility in University Timetabling

Thomas Feutrier<sup>1</sup>, Marie-Éléonore Kessaci<sup>1</sup>, Nadarajen Veerapen<sup>1</sup>

Univ. Lille, CNRS, Centrale Lille, UMR 9189 CRISTAL, F-59000 Lille, France  
`{thomas.feutrier,mkessaci,nadarajen.veerapen}@univ-lille.fr`

**Mots-clés** : *Timetabling, Feasibility Prediction, Instance Analysis, Metaheuristics*

## 1 Introduction

We consider the Curriculum-Based Course Timetabling (CB-CTT) problem, one of several NP-hard combinatorial problems found within University Timetabling. CB-CTT is based on the notion of a curriculum, which corresponds to a package of courses. Thus, students enroll in a curriculum as is the case in most French universities. The objective of the CB-CTT is to minimize the sum of its soft constraint violations.

The 2007 International Timetabling Competition (ITC 2007) featured a CB-CTT track and a set of benchmark instances that spurred research on the problem. However, ITC 2007 proposed only 21 (real-world) instances. This limited amount of data has curtailed the ability of the research community to create robust prediction models [2] that can be used in the context of automatic algorithm selection or configuration. In a recent paper, de Coster et al. [1] have expanded the set of real-world instances to 82. In addition, they introduced a new generator and proposed about 7000 generated instances.

Our preliminary work on these new instances has shown that a non trivial number of them, especially the synthetic ones, are not solvable by the construction heuristic of the winning approach of ITC 2007 [3], which is still considered state-of-the-art today [1]. Our current work focuses on predicting which instances are infeasible for our solver and identifying whether synthetic instances exhibit similar behavior to real-world ones. Since there are few real-world instances, the aim is to develop a predictive model that essentially relies on synthetic instances for training while still maintaining high accuracy when testing on real-world data. We find that using all the synthetic instances yields a poor model, indicating a mismatch between generated and real-world instances. However, we identify a subset of synthetic instances that produces a better model, and thus that more closely resembles the real-world instances.

## 2 Preliminary Analysis

We assess the ratio of feasible instances as a first step in comparing synthetic and real-world instances. Results show that 62% of the generated instances are feasible while this proportion rises to 93% for the real-world instances. This marked difference in the proportion of infeasible instances is an indication that the instance generator does not necessarily produce instances that match the behavior of real-world ones. To visualize the similarity between instances, we analyze their spatial distribution via Principle Component Analysis (PCA).

Figure 1 uses a PCA to transform the n-dimensional space of instance features (such as the number of teachers, number of curricula, ...) into two dimensions on which the instances are plotted. While there is an overlap between the two sets of instances, generated instances are clearly more spread out. To shed light on the disparities between these two sets, feasibility prediction models are created.

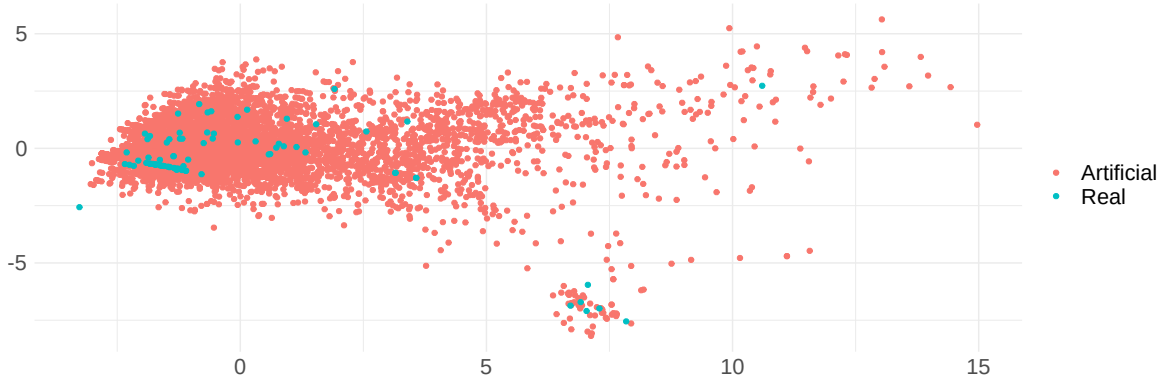


FIG. 1 – Visualization of Artificial and Real-world instances with a PCA on instance features.

### 3 Feasibility Prediction

We consider two models :  $\text{model}_{\text{Real}}$  and  $\text{model}_{\text{Artificial}}$ , which are respectively trained with only real-world and with only generated problem instances. Both are evaluated via 3-fold cross-validation.

The results of these models show that the accuracy, i.e. the fraction of correct predictions, is good, at about 0.78, when the test set uses the same type of instances as the train set. On the contrary, if  $\text{model}_{\text{Artificial}}$  is used to predict the feasibility of the real-world instances, then accuracy is very poor, or about 0.36. This behavior highlights a major difference between these two types of instances. Moreover, the analysis of the important features of both models shows that models do not favor the same features. The last part of this work focused on the analysis of the models' errors to improve them. Thus, the error distributions differed strongly from other instances. Therefore, in future work, we will surely be able to improve those models with recognizable errors.

### 4 Instance Selection

In this last part, the following hypothesis is investigated : some of the generated instances are probably similar quite similar to the real-world instances, and the poor generated instances can be filtered out of the initial set. Given the similarity of the feature distributions and the spatial proximity w.r.t. a PCA, a subset of the synthetic instances should be similar enough to real-world instances for training an efficient model that works on real-world instances. To identify those instances,  $\text{model}_{\text{Real}}$  is first used to predict the feasibility of artificial instances and 58 % of them are correctly predicted. The latter are used to train a new model, named  $\text{model}_{\text{Selected}}$ , which is used to predict the feasibility of real-world instances. Using this method, the accuracy improves from 0.36 to 0.77, with about 4000 selected training instances used. Since  $\text{model}_{\text{Selected}}$  works well, we set up a process that selects the right artificial instances more easily. The goal is to avoid the use of the  $\text{model}_{\text{Real}}$  and to use only the instance features.

## Références

- [1] A. De Coster, N. Musliu, A. Schaerf, J. Schoisswohl, and K. Smith-Miles. Algorithm selection and instance space analysis for curriculum-based course timetabling. 25(1) :35–58, 2022.
- [2] T. Feutrier, M.-E. Kessaci, and N. Veerapen. Analysis of search landscape samplers for solver performance prediction on a university timetabling problem. In *PPSN XVII*, volume 13398 of *LNCS*, pages 548–561. Springer International Publishing, 2022.
- [3] T. Müller. ITC2007 solver description : a hybrid approach. 172(1) :429, 2009.