

# Coursework Distributed Databases

Strijk Philippe, Wout Boeykens, Hannes Roegiers<sup>1</sup>

## Samenvatting

This task discusses 3 Kaggle competitions

## Sleutelwoorden

Apache — Spark — Maven — Big Data — Machine Learning — Kaggle — Distributed Databases

## Co-promotor

<sup>2</sup> (Van Vreckem Bert)

**Contact:** <sup>1</sup> wout.boeykens@student.hogent.be; <sup>2</sup> hannes.roegiers@student.hogent.be; <sup>3</sup> philippe.strijk@student.hogent.be

## Inhoudsopgave

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data Preprocessing</b>	<b>1</b>
<b>3</b>	<b>Model Selection</b>	<b>1</b>
3.1	Logistic Regression . . . . .	1
3.2	Random Forest Classifier . . . . .	1
<b>4</b>	<b>Conclusions</b>	<b>1</b>

## 1. Introduction

At the start of this task, our knowledge about Java Spark and the implementation of Machine Learning methods was quite limited, if not inexistant. Translating code from easy-to-use python scikit-learn to Java Apache Spark libraries proved to be a challenge that would send us down many knowledge-filled rabbit holes. The unique character of this assignment was evidently clear from the get-go. This coursework discusses 3 Kaggle competitions, namely the Quora Insincere Question Classification, <Wout Subject>, <Hannes Subject>.

## 2. Data Preprocessing

## 3. Model Selection

### 3.1 Logistic Regression

### 3.2 Random Forest Classifier

## 4. Conclusions