



# **Automated collision detection using machine learning algorithms on sensor data**

## **Project Report**

from the Course of Studies Angewandte Informatik  
at the Cooperative State University Baden-Württemberg Mannheim

by  
**Tim Schmidt**

August 2018

<b>Time of Project</b>	17 weeks
<b>Student ID, Course</b>	8531806, TINF15AI-BC
<b>Company</b>	IBM Deutschland GmbH, Ehningen
<b>Supervisor in the Company</b>	Julian Jung
<b>Reviewer</b>	Julian Jung

## Author's declaration

Hereby I solemnly declare:

1. that this Project Report, titled *Automated collision detection using machine learning algorithms on sensor data* is entirely the product of my own scholarly work, unless otherwise indicated in the text or references, or acknowledged below;
2. I have indicated the thoughts adopted directly or indirectly from other sources at the appropriate places within the document;
3. this Project Report has not been submitted either in whole or part, for a degree at this or any other university or institution;
4. I have not published this Project Report in the past;
5. the printed version is equivalent to the submitted electronic one.

I am aware that a dishonest declaration will entail legal consequences.

Mannheim, August 2018

---

Tim Schmidt

# Contents

<b>List of Figures</b>	<b>III</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Objective . . . . .	2
1.3 Limitations of the Research . . . . .	2
1.4 Organization of the Research . . . . .	2
<b>2 Theory</b>	<b>4</b>
2.1 Machine Learning . . . . .	4
2.2 Classification . . . . .	4
2.3 Evaluation of Classifiers . . . . .	4
2.4 Algorithms . . . . .	4
2.4.1 linear SVM . . . . .	4
2.4.2 logistic regression . . . . .	4
2.4.3 decision trees . . . . .	4
2.4.4 random forest . . . . .	4
2.4.5 gradient-boosted trees . . . . .	4
2.4.6 naive Bayes . . . . .	4
<b>3 Analysis of the pattern recognition problem</b>	<b>5</b>
3.1 Analysis of given sensor data . . . . .	5
3.2 Solution strategy . . . . .	5
<b>4 Algorithm configuration</b>	<b>7</b>
4.1 Domain based feature selection . . . . .	7
4.2 Parameter tuning . . . . .	7
<b>5 Evalutation</b>	<b>8</b>
<b>6 Outlook</b>	<b>9</b>

# List of Figures

# 1 Introduction

## 1.1 Motivation

The field of data analytics and machine learning is increasingly becoming the basis of industry competition. This is especially true in combination with another growing field, the Internet of Things (IoT), the McKinsey Global Institute's study 'The Age of Analytics: Competing In A Data-Driven World' found. [1] IBM has set a goal to tackle the challenges of IoT and Big Data with an array of new services, industry offerings and capabilities for enterprise clients, startups and developers. [2] These services are part of IBM Bluemix, a cloud platform as a service (PaaS).

An important step for IBM is getting clients interested and aware of the capabilities of IoT and machine learning so that collaborations and projects with clients can be acquired. One way this is done is by presenting showcases – small projects that show in an exemplary way the possibilities and values of the technology. It is essential to show technologies from the IBM Bluemix portfolio to present differentiating factors to competitors.

The Sensorboard is such a showcase. A physical skateboard is mounted with a Texas Instrument's SensorTag, and is able to connect wirelessly to the cloud platform IBM Bluemix. From there one or more Sensorboards can be monitored and managed. As part of the showcase a usage of the Sensorboard demonstrated as a rentable and cloud-managed vehicle and solution for urban mobility is.

This showcase will be extended with machine learning applications to present the value machine learning algorithms can have on sensor data. The concrete goal of that extension is to recognize collision that happen to the Sensorboard and that indicate accidents. This recognition will be based solely on sensor data, thus by recognizing patterns in the time series data. For this IBM Bluemix services are used to showcase the capabilities of the platform. The services which are requested to be demonstrated and used in the showcase are the IBM IoT Foundation, to connect the Sensorboard to the platform, and the IBM Data Science Experience as a development environment for an underlying Apache Spark cluster, on which machine learning models will be created.

## 1.2 Objective

A consumer grade skateboard is equipped with a Texas Instrument's SensorTag. The Texas Instrument's SensorTag collects acceleration data from the skateboard on three spacial axis. To process this data cloud services of IBM Bluemix are used. This especially includes the IBM Data Science Experience and an Apache Spark service. Consequently utilizing Apache Spark as a framework for cluster computing, machine learning models are trained using Apache Spark. This narrows down the selection of machine learning algorithms to the algorithms available in function library Apache Spark MLlib, which provides scalable algorithms, optimized to run on Apache Spark.

This paper investigates on the feasibility of machine learning algorithms for collision detection in the above described setting and identifies, if feasible, the best performing algorithm under the described conditions.

The scenario of impact detection that may indicate accidents requires a higher focus on detecting an accident than on preventing false alarms. A metric is found to represents this unequal relationship in the assessment of algorithm performance.

## 1.3 Limitations of the Research

The technologies described in the chapters Motivation and Objective like the Texas Instrument's SensorTag, the IBM Bluemix services and programming libraries are a given for the project by IBM. Therefore reasons for this particular decisions will not be discussed and no evaluation of alternatives will be conducted in this paper. The findings of this paper are acquired and only applicable for the specific hardware and configuration described in the paper. Findings can not be transferred to other vehicles, sensors, configurations or any other setup than the one described in this paper. The data which is used to come to decisions and to train models is generated in experiments, which are conducted solely for this reason. No other sensor data of skateboards or other vehicles are used in the considerations of this paper. The machine learning algorithms examined only include algorithms currently provided by the function library Apache Spark MLlib. A list of relevant algorithms currently provided can be seen in the appendix in chapter ??.

## 1.4 Organization of the Research

In the beginning the topic of machine learning is introduced and relevant machine learning algorithms and evaluators are explained with sufficient theoretical and mathematical

background information. Following, the pattern recognition problem will be analysed in detail. Bringing together the theoretical knowledge about machine learning and the outlined problem, a procedure is developed how the data is best prepared and algorithms are best configured and applied. The results of the application are then evaluated and compared to each other. The question of whether machine learning algorithms are feasible for collision detection and if so, which algorithm is best performing under the considered conditions, is then answered. In the end an outlook is given on possible further research.

# 2 Theory

## 2.1 Machine Learning

## 2.2 Classification

## 2.3 Evaluation of Classifiers

## 2.4 Algorithms

### 2.4.1 linear SVM

### 2.4.2 logistic regression

### 2.4.3 decision trees

### 2.4.4 random forest

### 2.4.5 gradient-boosted trees

### 2.4.6 naive Bayes



## 3 Analysis of the pattern recognition problem

### 3.1 Analysis of given sensor data

This chapter will provide a closer examination of the time series data belonging to the Sensorboard showcase. This data is basis for the further work of this paper and for answering the questions proposed in Chapter 1.1.

The sensor data measured by the Texas Instrument's SensorTag is collected and send to the IBM IoT Foundation. Measurements are taken at sampling rate of 10 measurements per second. To optimize accuracy the maximal sampling rate of the Texas Instrument's SensorTag is used[ref to TI technical docu]. The measurements include acceleration in spacial 3 axis. Those measurements are taken simultaneously and handled from then on together as one measurement collection. This collection is enriched with the time of the measurement with millisecond accuracy. A example measurement collection can be seen in ... ?? Using the provided time of a measurement the data can be plotted on a timeaxis as seen in ... ?. A measurement collection as it is described above will be also referred to as a data point in this paper. The single measurements will be also referred to as features of a data point.

The Sensorboard has been tested and sensor data has been collected. During the testing of the Sensorboard 69 collisions were simulated in total. Of the, in this way collected, 6453 datapoints 69 were labeled as a collision while the rest was labeled as no collision. The labeling was figured out by analysing video footage recorded during the test drive. The exact time of a collision in the video was matched with the time of the recorded sensor data to label data points correctly.

### 3.2 Solution strategy

Combining the knowledge about the sensor data from section 3.1 and the objective as described in chapter 1, this chapter will narrow down on a concrete problem type and solution strategy. In the process limitations described in Chapter 1 have to be considered. Having a set of prelabeled sensor data as described in section 3.1 allows for a supervised machine learning approach. Furthermore the existence of concrete data points which

represent the state of the Sensorboard at a given time make it possible to base the recognition of collisions solely on those data points. The recognition of incidents can then in turn be formulated as the problem of deciding whether a given data point corresponds to a collision or no collision. A problem of this kind can be solved by binary classification. Having specified the solution method to a binary classification model, concrete algorithms can be considered. In section ?? the necessity of an Apache Spark optimized algorithm and the usage of the Apache Spark library module MLlib is constituted. For the specific task of supervised binary classification the following methods are supported by the MLlib module:

1. Linear models
  - a) Linear SVM
  - b) Logistic regression
2. Decision tree
3. Ensemble of decision trees
  - a) Random forest
  - b) Gradient-boosted trees
4. Naive Bayes

Only the possible algorithms mentioned here will be parameterized and trained on the collected data in chapter ??.

To assess the performance of resulting models a clear metric will be used for comparison. This metric has to represent the unequal importance of detecting a collision and preventing false alarms. As mentioned in section ?? the act of detecting a collision that actually happened (true positives) is more important than preventing the detection of collision that didn't happen (false positives). In a classification setting the effectiveness in detecting collision that actually happened is described by recall. The effectiveness in preventing the detection of collision that didn't happen is described by precision. The metric that will be used to represent the unequal relationship between precision and recall is a  $f_2$ -score. This metric puts more weight on the recall of a model.

## **4 Algorithm configuration**

### **4.1 Domain based feature selection**

### **4.2 Parameter tuning**

## 5 Evalutation

To answer the question of feasability and the best performance the

## 6 Outlook