

# A Proposal to Mitigate Similarity Bias for the Paderborn Bearing Data Set

Lúcio Antônio Stange Venturim, Francisco de Assis Boldt  
Programa de Pós-graduação em Computação Aplicada (PPComp)  
Campus Serra do Instituto Federal do Espírito Santo (IFES)  
Rodovia ES-010 – Mangueiras – Serra – ES – Brazil  
lucioventurim@gmail.com, franciscoa@ifes.edu.br

**Abstract**—Similarity bias is a phenomenon that might occur when samples in a data set are originated from the same acquisition or are acquired from similar equipment settings. This phenomenon may lead to overoptimistic evaluations of machine learning algorithms. That overoptimistic estimation probably will not be reflected when the models generated by these algorithms will be put on production. This paper proposes a method to mitigate the similarity bias for bearing fault diagnosis evaluation with the Paderborn bearing dataset. The method consists on defining the proper data splits for training, validation and testing. Related works of classification models for the Paderborn data set were evaluated, showing that the similarity bias was not considered. Experiments were performed with K-nearest neighbors, Random Forest and *FaultNet* classification models. Accuracy and F1-score results show that the proposed method is effective to mitigate the similarity bias.

**Index Terms**—bearing, similarity bias, paderborn, fault diagnosis

## I. INTRODUCTION

Faults or unexpected changes might occur on equipment components during production, causing maintenance stops or even a breakdown [1]. Rolling bearings are indispensable for machines with rotating parts. Monitoring faults in these components is crucial, since up to 50% of failures on equipment with bearings are caused by faults on these parts [2]. Monitoring industrial processes can be divided into four stages: fault detection, fault identification, fault diagnosis and recovery process. Among these steps, fault diagnosis can be applied through software-based systems, which are considered essential tools to ensure the security and maintenance of dynamic processes [3]. These systems can use machine learning techniques, which use signals collected from equipment to train and test classifiers [4].

A typical data for bearing damage detection is the vibration signal acquired through acceleration sensors. The data may be composed by several signals from different parts. Depending on how the signals were extracted, they must be split in samples, in order to have enough examples to train, validate and test a classification algorithm. However, depending on the data splitting strategy, the experiments might result in overoptimistic evaluations. For these cases, the results would be suitable only for specific conditions. Thus, when the developed model is applied in practice it will not have a performance as good as that obtained during the training and testing phases.

A realistic assessment must also take into account the generalization capacity of the developed model. The generalization refers to a model's ability to adapt properly to new, previously unseen data, drawn from the same distribution as the one used to create the model. The trained classifier must be able to recognize faults in as many conditions as possible, even when there are variations in their working conditions. A characteristic that influences the generalization ability of an algorithm, according to [4], is the similarity bias. Similarity bias occurs when data from the same condition are used both for training and testing the model and have very similar characteristics, making the classification task relatively trivial. This article proposes a method to mitigate the similarity bias for a data set publicly available by the University of Paderborn.

Section II explains the concept of similarity bias. Section III presents details of the Paderborn data set and related published works. Section IV describes the proposed method to mitigate the similarity bias for the Paderborn data set. Section V presents the results and its analysis, with comparisons to the results of other works. Finally, section VI concludes the article and discusses possibilities for future works.

## II. THE SIMILARITY BIAS

In machine learning, bias is the phenomenon of observing results that are systematically prejudiced due to erroneous assumptions in the learning process. Several types of bias may occur, in many shapes and forms. Examples are historical bias, representation bias, measurement bias, among others [5]. This work focuses on the similarity bias, which was defined by [4] as a phenomenon that might occur when samples in a data set are originated from the same acquisition or are acquired from similar equipment settings. Samples representing the same condition and acquired from the same acquisition are typically very similar. This may lead to overoptimistic evaluations of machine learning algorithms. That overoptimistic estimation probably will not be reflected when the models generated by these algorithms are put on production.

In [4], experiments were performed with a data set provided by the Case Western Reserve University (CWRU), with different data splits, in order to check the occurrence of similarity bias. The similarity bias in the CWRU data set may occur, according to [4], because the samples carry information of the acquisition itself, and not only about the failure. The