
Use of machine learning to study batch effect removal methods in metabolomics data

Marc Solé Estragues, 202002465

Master's Thesis, Bioinformatics

Supervisor: Palle Villesen

June 2022



AARHUS
UNIVERSITET
CENTER FOR BIOINFORMATIK (BiRC)

Abstract

The fusion of criminology and machine learning is a fascinating field, the use of machine learning models able to identify and predict biological characteristics is changing the way criminal investigations are handled. Using blood metabolomics data it is possible to construct machine learning models able to interpret the chemical changes that happen as the sample degrades through time and are able to predict the age of a blood spot. This information can then be used to determine the time when a crime took place. One of the big limitations of metabolomics data is the presence of batch effect. The removal of such undesired effect is a necessity in order to obtain unbiased data which can later be used in comparative analyses which require of high statistical power such as the construction of machine learning models. Many different techniques to achieve this correction on the data are available, and which one should be used to obtain the best data capable of the best machine learning performance is a question that needs to be answered.

This thesis has sought to study how these methods can be applied and how they transform a specific set of blood metabolomics data obtained through HPLC-MS/MS which is being used in research to construct a machine learning model able to predict the age of dried blood spots. Of special interest has been to compare how well different methods are able to remove the batch effect from the samples, as well as how much improvement is observed in the predictive error using simple elastic-net and random forests regression models.

It was determined through principal component analysis that *ComBat* and *WaveICA2.0*, two of the most advanced batch removal methods available, as well as mean centering, perform good batch correction on the data, with other methods such as row normalisation and probabilistic quotient normalisation performing the worst in that regard. The analysis also revealed that when applying *WaveICA2.0* the correct cutoff parameter needs to be chosen, and much work went into exploring how these cutoff value changes the amount of batch correction applied. In the end the PCAs revealed that the data does contain batch effect but it is orthogonal to the biological variable of age being studied, and thus when the batch correction is applied little change in the statistical performance is observed and the root mean square error obtained from the ML model's predictions can even increase in some cases.

Over-all, the results show that some batch correction methods perform better than others, but their use will depend on the final aim of the project and the type of data that is being worked with. In the case of this specific metabolomics data, strong batch correction results in lower predictive performance overall and methods such as *WaveICA2.0* which have a tuning capability in the cutoff parameter could be considered to adjust for the best amount of batch correction while maintaining good prediction accuracy.

Aknowledgements

I dedicate this work of non fiction to my family who from the distance have always supported me, my bros Adam Petro and Lisa Berling who have always motivated me to keep going, to my supervisor Palle Villesen for his great guidance and advice, and finally I want to thank the whole country of Denmark and Aarhus University for welcoming me with open arms.

*Marc Solé Estragués,
Aarhus, 12th June 2022.*

Contents

1	Introduction	1
1.1	Introduction to metabolomics	1
1.2	Brief description of the Mass Spectrometry analytic technique	2
1.3	Metabolomics analyses of blood spots in the criminology field	4
1.4	Introduction to Batch Effects	5
1.5	Current Batch Effects removal methods	7
1.6	Objective	8
2	Materials and Methods	10
2.1	Data	10
2.2	General procedure	11
2.3	Basic quality control	12
2.4	Types of batch correction techniques used	13
2.4.1	Row normalisation	13
2.4.2	Quantile normalisation	13
2.4.3	Mean normalisation	13
2.4.4	<i>ComBat</i>	14
2.4.5	<i>WaveICA2.0</i>	15
2.4.6	Probabilistic Quotient Normalisation	16
2.5	Unsupervised approaches to visualize batch effect	17
2.5.1	Principal Component Analysis	17
2.5.2	Metabolite batch visualization	18
2.6	Statistical modelling	18
2.6.1	The <i>caret</i> package	19
2.6.2	Models used	19
2.6.3	Procedure	20
2.7	Additional material: <i>WorkflowR</i>	22
3	Results and Discussion	23
3.1	General overview of the data	23
3.1.1	Principal component analysis on <i>Batch</i> and <i>Age</i>	25
3.1.2	Feature intensity analysis by injection order and batch	25
3.1.3	Machine Learning performance	26
3.2	Results after batch correction	27
3.2.1	Principal Component Analysis	27
3.2.2	Feature Analysis	33
3.2.3	Model performance	38
3.2.4	Analysing how the age variable changes after batch correction	44
4	Final conclusions	46
A	Additional Material	54