# Bayesian Non-parametric Models for Time Segmentation and Regression

## Ava Bargi

MEng    Baha'i Institute for Higher Education (BIHE)

Faculty of Engineering and Information Technology

University of Technology, Sydney

*A thesis submitted in fulfilment of the requirements for the degree of*

***Doctor of Philosophy***

November 2015

*To **Nima**,*

*a wonderful companion,*

*a patient and gentle support,*

*an unconditional lover,*

*my husband and most precious friend,*

*to whom I am willing to dedicate this thesis and all my future*

*accomplishments !*

# Certificate of Authorship and Originality

Title: **Bayesian Non-parametric Modelling for Infinite-modal Segmentation and Prediction**

Author: **Ava Bargi**

Date: **November 27, 2015**

Degree: **PhD**

I certify that the work in this dissertation has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the dissertation has been written by me. Any help that I have received in my research work and the preparation of the dissertation itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the dissertation.

Signiture of author

# Acknowledgements

I would like to sincerely thank the following wonderful people, without whose attention and support this thesis could not have been completed:

My principal supervisor, *Professor Massimo Piccardi*, for gently holding my hands in the toddling stages of a PhD; for constantly caring for my progress and happiness; for teaching me the real essence of research, far beyond publish or perish: the never changing virtue of *'modest learning'*. In addition to a smart supervisor, he has been a great friend and company. It has been both a joy and honour to be his student and I am utterly grateful for all the learning during this PhD.

*Dr Richard Xu*, my co-supervisor, has been a wonderful technical adviser and support. His timely advice has impacted the most valuable accomplishments. He has been another role model for genuine spirit of research and hard work.

*Dr Emily B. Fox* whose research has largely inspired this work, offering me an example of noble ideas, elegant presentation and helpful and friendly attitude.

*Professor Zoubin Ghahramani* for organising a talk in his center, offering me an opportunity to communicate the current research and receive invaluable feedback. His modest and open attitude to collaboration has marked an important point in my research career.

*Professor Yee Whye Teh, Dr Christoph Lampert, Dr Fernando De la Torre*

*Dr Nima Khorsandnia*, who has always encouraged and reassured me in the critical moments. I feel really prosperous to have his company in life.


Ava Bargi,

March 2015, Sydney

# Abstract

Non-parametric Bayesian modelling offers a principled way for avoiding model selection such as pre-defining the number of modes in a mixture model or the optimal number of factors in factor analysis. Instead, Bayesian non-parametric methods allow the data to determine the complexity of model. In particular, the hierarchical Dirichlet process (HDP) is used in a variety of applications to infer an arbitrary number of classes from a set of samples. Within the temporal modelling paradigm, Bayesian non-parametrics is used to model sequential data by integrating HDP priors into state-space models such as HMM, constructing HDP-HMM. Also in latent factor modelling and dimensionality reduction, Indian buffet process (IBP) is a well-known method capable of sparse modelling and selecting an arbitrary number of factors among the often high-dimensional features.

In this PhD thesis, we have applied the above methods to propose novel solutions to two prominent problems. The first model, named as 'ADON HDP-HMM', is an *adaptive online* system based on HDP-HMM. 'ADON HDP-HMM' is capable of segmenting and classifying the sequential data over unlimited number of classes, while meeting the memory and delay constraints of streaming contexts. The model is further enhanced by a number of *learning rate*s, responsible for tuning the adaptability by determining the extent to which the model sustains its previous parameters or adapts to the new data. Empirical results on several variants of synthetic

and action recognition data, show remarkable performance, particularly using adaptive learning rates for evolutionary sequences.

The second proposed solution is an elaborate factor regression model, named as *non-parametric conditional factor regression* (NCFR), to cater for multivariate prediction, preserving the correlations in the response layer. NCFR enhances factor regression by integrating IBP to infer the optimal number of latent factors, in a sparse model. Thanks to this data-driven approach, NCFR can significantly avoid over-fitting even in cases where the ratio between the number of available samples and dimensions is very low. Experimental results on three diverse datasets give evidence of its remarkable predictive performance, resilience to over-fitting, good mixing and computational efficiency.

# Contents

# List of Figures

# List of Tables