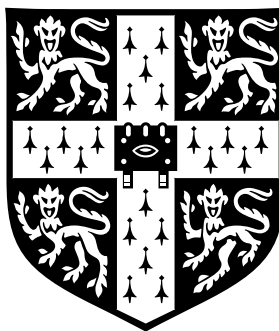


# Probabilistic Document Modelling

Philip J. Cowans  
Churchill College  
Cambridge

A dissertation submitted in candidature for the degree of Doctor of Philosophy,  
University of Cambridge

Inference Group  
Cavendish Laboratory  
University of Cambridge



January 2006



# DECLARATION

I hereby declare that my dissertation entitled “Probabilistic Document Modelling” is not substantially the same as any that I have submitted for a degree or diploma or other qualification at any other University.

I further state that no part of my dissertation has already been or is being concurrently submitted for any such degree or diploma or other qualification.

Except where explicit reference is made to the work of others, this dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration. This dissertation does not exceed sixty thousand words in length.

Date: .....

Signed: .....

Philip J. Cowans  
Churchill College  
Cambridge  
January, 2006

# ABSTRACT

In this thesis the development and application of probabilistic models of documents is considered. The initial focus is on language models which provide a way of modelling plain text documents. In particular the hierarchical Dirichlet language model, which is derived from simple Bayesian theory, is investigated and is shown to be well approximated by an existing method known as generalised PPM-A. Using this equivalence, generalised PPM-A is extended to produce a language model which while working on the level of individual letter-like symbols is able to make use of the division of the text stream into words. It is shown that the new model can be used in conjunction with a word list to improve performance when very little information from which to learn the statistics of the language is available.

The hierarchical Dirichlet model is then applied to the task of information retrieval, producing a new retrieval method which naturally includes document frequency information. This information has traditionally been used in retrieval systems, but previously had either been missing or introduced heuristically in language model based approaches to the problem. The hierarchical approach is also extended to the task of retrieval at the passage level where it is shown to give promising results.

Finally, the scope of the investigation is broadened to include documents which contain diagrams as well as plain text. A method is developed to group fragments of digitised ink strokes into perceptually relevant components of a diagram, while at the same time labelling the components with an object class. The approach, which is based on the conditional random field, is shown to work well both in terms of grouping and improving labelling performance when compared to other methods.

# ACKNOWLEDGEMENTS

I would like to thank Ryan Adams, Tom Minka, Iain Murray, Edward Ratzer, Oliver Stegle, David Stern, Martin Szummer, Keith Vertanen, Hanna Wallach, David Ward, Seb Wills and John Winn for ideas, suggestions, discussions and feedback without which the work presented in this thesis would not have been possible.

During my studies I was lucky enough to be able to spend three months working in the Machine Learning and Perception group at Microsoft Research UK, an opportunity for which I am grateful and which produced many of the ideas presented in Chapter 6.

Above all, I would like to thank my supervisor, Professor David MacKay for an unending supply of inspiration, discussion and feedback, and for providing invaluable advice on the draft of this thesis. I would also like to thank my parents, Sue and Trevor Cowans, and my brother, Nick, for their support and encouragement during the last four years.

This work was supported by a grant from Microsoft Corporation.

# CONTENTS

<b>Chapter 1</b>	<b>Introduction</b>	<b>1</b>
1.1	Why probabilistic models? . . . . .	1
1.2	Fundamentals of probability theory . . . . .	2
1.3	Applications of document modelling . . . . .	3
1.3.1	Language modelling . . . . .	3
1.3.2	Multimedia documents . . . . .	8
<b>Chapter 2</b>	<b>Language Modelling</b>	<b>9</b>
2.1	Introduction . . . . .	9
2.2	Measuring language model performance . . . . .	10
2.3	Review of existing techniques . . . . .	11
2.3.1	Maximum likelihood estimation . . . . .	11
2.3.2	Simple Bayesian approaches . . . . .	13
2.3.3	Smoothing . . . . .	15
2.3.4	Linear interpolation . . . . .	19
2.3.5	PPM and Witten–Bell smoothing . . . . .	19
2.3.6	Katz and Church–Gale smoothing . . . . .	20
2.3.7	Absolute discounting and Kneser–Ney smoothing . . . . .	21
2.3.8	Reduced context dimensionality . . . . .	21
2.3.9	Maximum entropy language models . . . . .	22
2.3.10	Topic models . . . . .	23
2.4	Experimental evaluation of generalised PPM-A . . . . .	23
2.4.1	Baseline performance and varying alpha . . . . .	24
2.4.2	Evaluation of update exclusion . . . . .	26
2.5	Hierarchical Bayesian methods . . . . .	27
2.5.1	The hierarchical Dirichlet distribution . . . . .	27
2.5.2	Deeper hierarchies . . . . .	28
2.5.3	The hierarchical Dirichlet language model . . . . .	29
2.5.4	Pólya urns and oracles . . . . .	30
2.5.5	Sampling from the hierarchical Dirichlet distribution . . . . .	31
2.6	Bayesian interpretation of generalised PPM-A . . . . .	32

---

2.6.1	Backing off and interpolation . . . . .	34
2.7	Empirical evaluation . . . . .	34
2.7.1	Comparison of the two approximations . . . . .	35
2.7.2	Investigating the prior distribution . . . . .	40
2.7.3	Posterior distribution . . . . .	42
2.7.4	Sampling results . . . . .	46
2.8	A note on the two parameter model . . . . .	49
2.9	Conclusions . . . . .	49
<b>Chapter 3</b>	<b>Word-Based Language Modelling</b>	<b>51</b>
3.1	Introduction . . . . .	51
3.2	Theoretical development . . . . .	51
3.2.1	Infinite Dirichlet models . . . . .	52
3.2.2	Top level priors . . . . .	53
3.2.3	Using a word list . . . . .	54
3.2.4	Update exclusion . . . . .	54
3.3	Implementation details . . . . .	55
3.4	Results . . . . .	56
3.5	Conclusions . . . . .	58
<b>Chapter 4</b>	<b>Text Retrieval</b>	<b>60</b>
4.1	Introduction . . . . .	60
4.2	Review of previous work . . . . .	60
4.2.1	Vector space models . . . . .	60
4.2.2	Binary independence retrieval . . . . .	61
4.2.3	Language modelling approaches . . . . .	63
4.2.4	Relevant sets . . . . .	64
4.2.5	The probability ranking principle . . . . .	64
4.2.6	Relevance feedback . . . . .	65
4.3	Whole collection models . . . . .	65
4.3.1	The model . . . . .	67
4.3.2	Experimental evaluation . . . . .	68
4.4	Results . . . . .	69
4.4.1	Precision experiments . . . . .	69
4.4.2	Effect of varying $\alpha_1$ . . . . .	72
4.4.3	Precision–recall curves . . . . .	72
4.4.4	Variation of query length . . . . .	73
4.5	Deeper hierarchies and passage retrieval . . . . .	75
4.5.1	Query term weighting . . . . .	77
4.5.2	Document length normalisation . . . . .	77
4.5.3	Implementation . . . . .	78

---

4.5.4	Whole document retrieval . . . . .	78
4.5.5	Passage selection . . . . .	80
4.6	Application to a seminar database . . . . .	83
4.7	Conclusions . . . . .	84
<b>Chapter 5</b>	<b>Dasher As A Search Tool</b>	<b>85</b>
5.1	Introduction . . . . .	85
5.2	Interactive search . . . . .	85
5.2.1	Data structure . . . . .	87
5.2.2	Storage requirements . . . . .	88
5.2.3	Dynamic node allocation . . . . .	92
5.2.4	Distributions over strings . . . . .	93
5.3	Contact list lookup . . . . .	94
5.3.1	Model adaption . . . . .	95
5.4	Further work . . . . .	96
5.5	Conclusions . . . . .	98
<b>Chapter 6</b>	<b>Electronic Ink Analysis</b>	<b>99</b>
6.1	Introduction . . . . .	99
6.2	Overview of the algorithm . . . . .	101
6.3	Undirected graphical models . . . . .	101
6.3.1	Conditional random fields . . . . .	103
6.3.2	Models over partitions . . . . .	104
6.4	Theoretical development . . . . .	105
6.4.1	Training . . . . .	107
6.4.2	Inference . . . . .	108
6.5	Operations over labelled partitions . . . . .	109
6.5.1	Message passing . . . . .	109
6.5.2	Message passing for labelled partitions . . . . .	111
6.5.3	Sum-product algorithm . . . . .	113
6.5.4	Message passing . . . . .	114
6.5.5	Max-product algorithm . . . . .	114
6.6	Complexity and the edge-dual representation . . . . .	115
6.6.1	Complexity . . . . .	117
6.7	Implementation details . . . . .	117
6.7.1	Graph construction . . . . .	118
6.7.2	Feature set . . . . .	118
6.7.3	Priors . . . . .	119
6.8	Experimental evaluation . . . . .	120
6.8.1	Learned weights . . . . .	120
6.8.2	Error rates . . . . .	121



---

6.9	Discussion . . . . .	122
6.10	Future extensions . . . . .	122
6.10.1	Additional features . . . . .	122
6.10.2	Tree-width constraints . . . . .	122
6.10.3	Approximate inference . . . . .	123
6.11	Conclusions . . . . .	124
<b>Chapter 7</b>	<b>Conclusions</b>	<b>129</b>
<b>Appendix A</b>	<b>Data Sets</b>	<b>131</b>
A.1	Language modelling . . . . .	131
A.1.1	Enron e-mail corpus . . . . .	131
A.1.2	Canterbury corpus . . . . .	131
A.2	Information retrieval . . . . .	132
A.2.1	TREC corpus . . . . .	132
A.2.2	Cranfield corpus . . . . .	132
<b>Appendix B</b>	<b>Stop Word List</b>	<b>133</b>
	<b>Bibliography</b>	<b>135</b>