

Module	<b>4F13</b>	Title of report	<b>Latent Dirichlet Allocation</b>			
Date submitted: <b>04/12/2020</b>		Assessment for this module is <input checked="" type="checkbox"/> <b>100%</b> / <input type="checkbox"/> 25% coursework of which this assignment forms <u><b>33</b></u> %				
<b>UNDERGRADUATE STUDENTS ONLY</b>		<b>POST GRADUATE STUDENTS ONLY</b>				
Candidate number:	<b>5562E</b>	Name:			College:	

Feedback to the student

☐ See also comments in the text

		Very good	Good	Needs improvmt
C O N T E N T	<b>Completeness, quantity of content:</b> Has the report covered all aspects of the lab? Has the analysis been carried out thoroughly?			
	<b>Correctness, quality of content</b> Is the data correct? Is the analysis of the data correct? Are the conclusions correct?			
	<b>Depth of understanding, quality of discussion</b> Does the report show a good technical understanding? Have all the relevant conclusions been drawn?			
	Comments:			
P R E S E N T A T I O N	<b>Attention to detail, typesetting and typographical errors</b> Is the report free of typographical errors? Are the figures/tables/references presented professionally?			
	Comments:			

Overall assessment (circle grade)	A*	A	B	C	D
Guideline standard	>75%	<b>65-75%</b>	55-65%	40-55%	<40%
Penalty for lateness:		20% of marks per week or part week that the work is late.			

Marker:

Date:

# 4F13 Probabilistic Machine Learning - Latent Dirichlet Allocation

Candidate: 5562E

November 26, 2020

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Questions</b>	<b>1</b>
a	Maximum Likelihood . . . . .	1
b	Bayesian Inference . . . . .	1

## 1 Introduction

We have a document test set  $\mathcal{A}$ , which consists of  $D$  documents indexed by  $d \in \{1 \dots D\}$ . A document is an ordered list of words. The vocabulary  $\mathcal{M}$  has  $M = |\mathcal{M}|$  unique words. We denote the  $n$ 'th word in the  $d$ 'th document by  $w_{nd} \in \{1 \dots N_d\}$ . Where  $N_d$  is the length of document  $d$ . For simplicity we denote the count of occurrences of word  $m$  in the test set by  $c_m$ .

We hold back a test set  $\mathcal{B}$  to calculate the performance of our approaches.

## 2 Questions

### a Maximum Likelihood

We begin by assuming that each word is drawn independently from a categorical distribution with parameter  $\beta$ :  $w_{nd} \stackrel{iid}{\sim} \text{Cat}(\beta)$ . In this case  $\beta$  is a  $M \times 1$  vector with the conditions that  $\sum_{m=1}^M \beta_m = 1$  and  $\beta_i \geq 0$ . The likelihood of the parameter  $\beta$  is the probability of the dataset given  $\beta$ :

$$L(\beta) = P(\mathcal{A}|\beta) = \prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{nd}|\beta) = \prod_{m \in \mathcal{M}} \beta_m^{c_m} \quad (1)$$

Where  $c_m$  is the count of word  $m$  in the training set. We wish to obtain the Maximum-Likelihood estimate  $\hat{\beta}^{ML} = \arg \max L(\beta)$ . We prefer to maximise the log-likelihood as this is more tractable:

$$\mathcal{L}(\beta) = \log L(\beta) = \sum_{m \in \mathcal{M}} c_m \log \beta_m \quad (2)$$

We can now take derivatives and include a Lagrange multiplier to respect the sum to 1 constraint:

$$\begin{aligned} \frac{\partial}{\partial \beta_i} \left\{ \mathcal{L}(\beta) + \lambda \left( 1 - \sum_{m=1}^M \beta_m \right) \right\}_{\beta=\hat{\beta}^{ML}} &= \frac{c_i}{\hat{\beta}_i^{ML}} - \lambda = 0 \\ \therefore \hat{\beta}_i^{ML} &= \frac{c_i}{\lambda} = \frac{c_i}{\sum_{m \in \mathcal{M}} c_m} = \frac{c_i}{C} \end{aligned} \quad (3)$$

Therefore, the ML estimate is simply the empirical frequency of each word (normalised by the sum of all counts  $C$ ).

### b Bayesian Inference

Words: XX

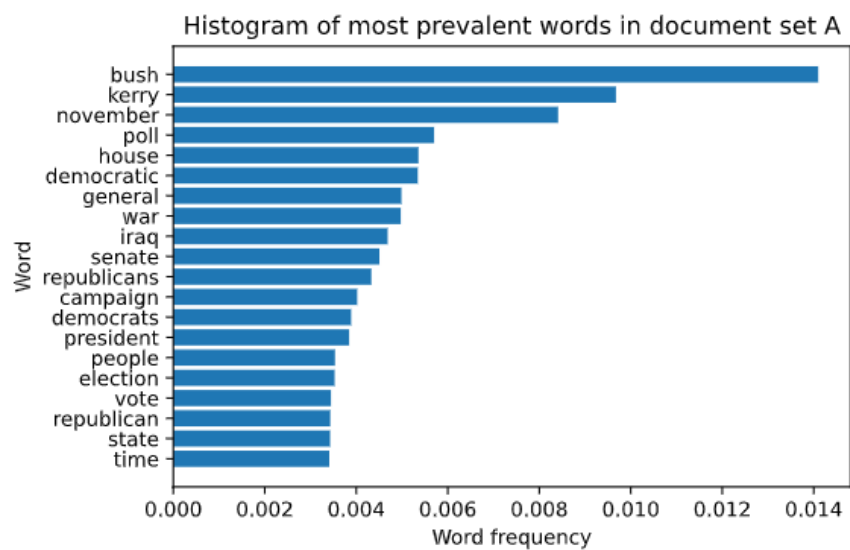


Figure 1: Histogram of top 20 most prevalent words in test set A