

USING TOPIC MODELS TO DESCRIBE DISTURBANCE & QUANTIFY RESPONSES TO ENVIRONMENTAL CHANGE

Gavin L. Simpson & Emma Wiik · University of Regina
ESA Portland 2017 · August 9th 2017

ACKNOWLEDGEMENTS



**NSERC
CRSNG**

Slides: bit.ly/esatopicmodels

Copyright © (2017) Gavin L. Simpson Some Rights Reserved

Unless indicated otherwise, this slide deck is licensed under a Creative Commons Attribution 4.0 International License.



COMMUNITY RESPONSE TO ENVIRONMENTAL CHANGE



By Mauricio Antón [CC BY 2.5 (<http://creativecommons.org/licenses/by/2.5>)], via Wikimedia Commons

COMMUNITY RESPONSE TO ENVIRONMENTAL CHANGE



By Credit: NOAA George E. Marsh Album [Public domain], via Wikimedia Commons

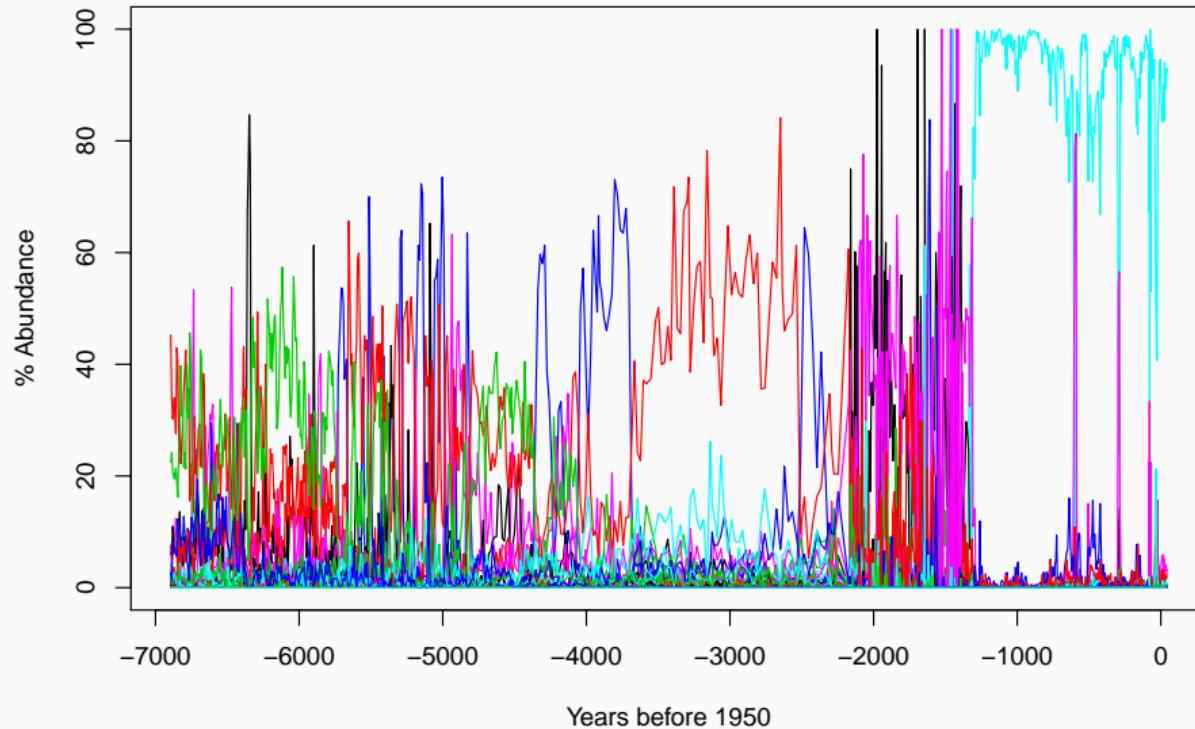
COMMUNITY RESPONSE TO ENVIRONMENTAL CHANGE



By Michael Knall [Public domain], via Wikimedia Commons

How have aquatic communities responded to these rapid changes over the last few millennial?

COMPLEX MULTIVARIATE SPECIES DATA



Data provided by Jeffery Stone (Indiana State University)

DIMENSION REDUCTION

Typically we can't model all 100+ taxa in data sets like this

- (M)ARSS-like models don't like the large n
- Most methods assume regular spacing in time

Seek a reduced dimensionality of the data that preserves the signal

Existing dimension reduction methods aren't appropriate for questions we want to ask

- Interpretation of latent factors is complex (PCA, CA, Principal Curves)
- Impose hard thresholding on the data (hierarchical clustering)
- No grouping of species or samples, or group only one

Can we group species into J **associations** and soft cluster sample as compositions of these associations?

Foy Lake

- Deep, freshwater lake
- Drought-sensitive Flathead River Basin
- Diatom assemblages sensitive to lake depth variation
- Related to variability in effective moisture

Regime shift ~1.3ka BP

Spanbauer *et al* PLOS ONE 9(10) e108936

OPEN ACCESS Freely available online

PLOS ONE



Prolonged Instability Prior to a Regime Shift

Trisha L. Spanbauer^{1*}, Craig R. Allen², David G. Angeler³, Tarsha Eason⁴, Sherilyn C. Fritz⁵, Ahjond S. Germestani⁶, Kirsty L. Nash⁵, Jeffery R. Stone⁶

¹ Department of Earth and Atmospheric Sciences, University of Nebraska-Lincoln, Lincoln, Nebraska, United States of America, ² U.S. Geological Survey, Nebraska Cooperative Fish and Wildlife Research Unit, University of Natural Resources, University of Nebraska-Lincoln, Lincoln, Nebraska, United States of America, ³ Department of Aquatic Sciences and Assessment, Swedish University of Agricultural Sciences, Uppsala, Sweden, ⁴ Office of Research and Development, National Risk Management Research Laboratory, U.S. Environmental Protection Agency, Cincinnati, Ohio, United States of America, ⁵ Australian Research Council Centre of Excellence for Coral Reef Studies, James Cook University, Townsville, Queensland, Australia, ⁶ Department of Earth and Environmental Systems, Indiana Tech University, Fort Wayne, Indiana, United States of America

Abstract

Regime shifts are generally defined as the point of ‘abrupt’ change in the state of a system. However, a seemingly abrupt transition can be the product of a system reorganization that has been ongoing much longer than is evident in a statistical analysis. We used a dataset of diatom assemblage abundance and environmental variables to examine the onset of a long-term high-resolution paleoecological dataset with a known onset of regime shift. Analysis of this dataset with Fisher information and multivariate time series modeling showed that there was a ~200-yr period of instability prior to the regime shift. This period of instability and the subsequent regime shift coincided with regional climate change and the onset of the Little Ice Age. This study extends the use of paleolimnology as a unique approach to test tools for the detection of thresholds and stable states, and thus to examine the long-term stability of systems over periods of multiple millennia.

Citation: Spanbauer TL, Allen CR, Angeler DG, Eason T, Fritz SC, et al. (2014) Prolonged Instability Prior to a Regime Shift. PLOS ONE 9(10): e108936. doi:10.1371/journal.pone.0108936

Editor: John A. D. Arora, University of Cambridge, United Kingdom

Received: July 11, 2014 Accepted: September 5, 2014 Published: October 1, 2014

This is an open-access article. All of its content is freely available and may be reprinted, distributed, modified, built upon, or otherwise used by anyone for any lawful purpose. The work is made available under the Creative Commons CC-BY public domain dedication.

Funding: The Nebraska Cooperative Fish and Wildlife Research Unit is jointly supported by a cooperative agreement between the U.S. Geological Survey, the University of Nebraska-Lincoln, and the State of Nebraska. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. This work was supported in part by the August T. Larson Foundation of the Institute University of Agricultural Sciences and the NFWF Integrative Inland Lakes Education and Research Translational (IERT) program (NFWF-0400046) and the Sedimentary Geology and Palaeolimnology program (NFWF-0123-1678). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

* email: tspanbau@unl.edu

Introduction

Ecosystems can undergo regime shifts and reorganize into an alternative state when a critical threshold is exceeded [1–3]. Most quantitative regime shift research has focused on abrupt shifts that have occurred during a period of human influence, but this has resulted in a limited number of new fast variable (e.g., nutrient loading, crude oil release), but it hasn't addressed how slow variables (e.g., long-term changes in climate) can alter ecosystem state. Palaeolimnological records can provide insight on the frequency and duration of transitions between alternative states in systems that are slow to both fast and slow variables, at timescales not accessible in the observed record.

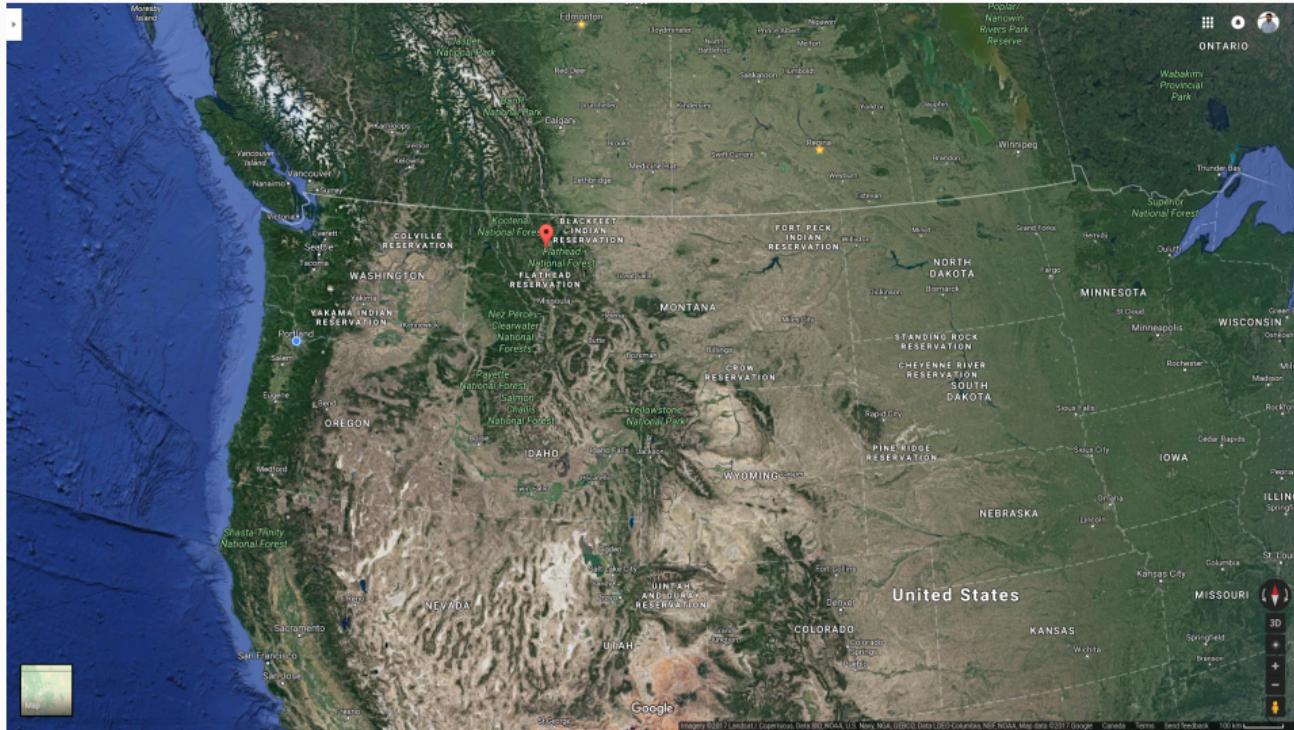
To test for regime shifts in the paleolimnological record, we used a long-term high-resolution sedimentological record from Foy Lake (Montana, USA) that showed abrupt changes in diatom community structure at ~1.3 ka (about 400 years before present, with present = AD 2010) [4]. The Lake Foy record (Fig 1; N. 41° 45' W, 105° 00' E) is a deep freshwater lake situated in the drought-sensitive Flathead River Basin in the Northern Rocky Mountains [4,5]. Diatom assemblages in this system are sensitive to changes in lake depth driven by changes in effective moisture [6] and represent one metric of ecological resilience. The percent abundances of 109 diatom species were collected from a lake sediment core that was sampled continuously at an interval of roughly ~5–20 years, yielding a ~1-ky record of 800 time-segments.

The paleolimnological data set we (i) plotted several indices proposed to be early-warning signals of approaching critical thresholds (‘warning variance’, ‘dashed responses’, ‘hysteresis, and the autocorrelation at lag-1’) [7] against time, (ii) collapsed the data species into three groups (freshwater, brackish, and saline), and

(iii) used multivariate time series modeling based on canonical ordination [8]. Many of these statistical early-warning signals have been developed based on bifurcation theory, and they have successfully anticipated regime shifts in nature [10–13], but not all [14].

Time-series data can be useful to infer the timing of the regime shift because changes between two different states prior to a regime shift [15]. Along with autocorrelation at lag-1, increasing variance in time-series data can be caused by critical slowing-down, where a system is slow to recover from minor disturbances and it approaches a steady state [7]. The time-series data can be useful to infer the stability, because critical signals often occur at the onset of the regime shift, which is generally too late to implement effective management actions [16]. Hence, we sought methods [7] and multivariate time series modeling that more effectively investigate the dynamics of complex multivariate systems. FL is an integral index based on information theory, declines as it approaches a

FOY LAKE — MONTANA



FOY LAKE — MONTANA



Machine learning approach for organizing text documents

Latent Dirichlet Allocation (LDA) – (Blei, Ng, & Jordan, *J. Mach. Learn. Res.* 2003)

generative model

Valle, D., Baiser, B., Woodall, C. W. & Chazdon, R. (2014) *Ecology Letters* 17

COMMUNITY OF SKITTLES



INDIVIDUAL SKITTLES FROM ONE OF FOUR FLAVOURED PACKS



What are the proportion of flavours in each pack?

How many of each pack comprise the skittle community?

LATENT DIRICHLET ALLOCATION

Aim is to infer the

- distribution of skittle **flavours** within each pack β_j , and
- distribution of skittle **packs** within each community (*sample*)

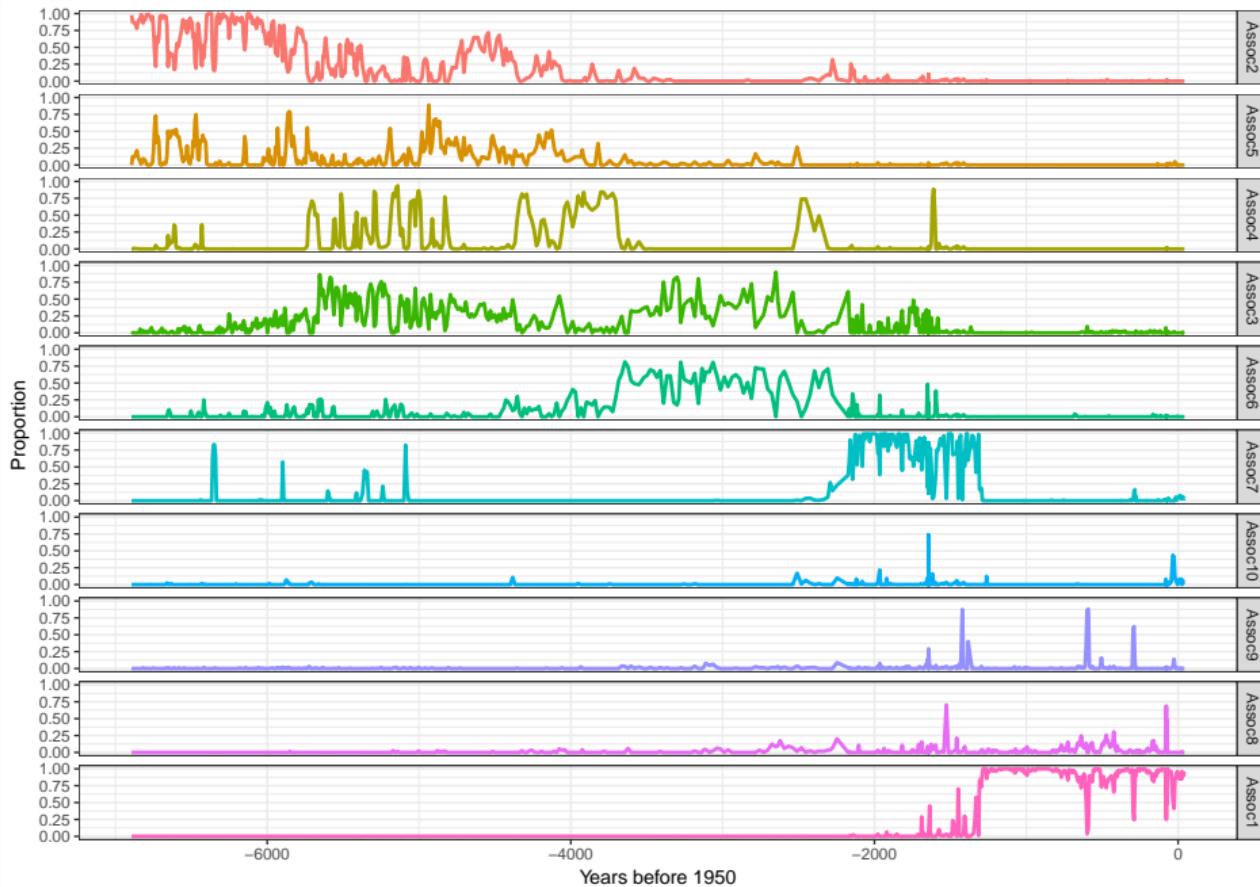
Achieve a soft clustering of samples – mixed membership model

Achieve a soft clustering of **species** (flavours) into **associations** of taxa (packs)

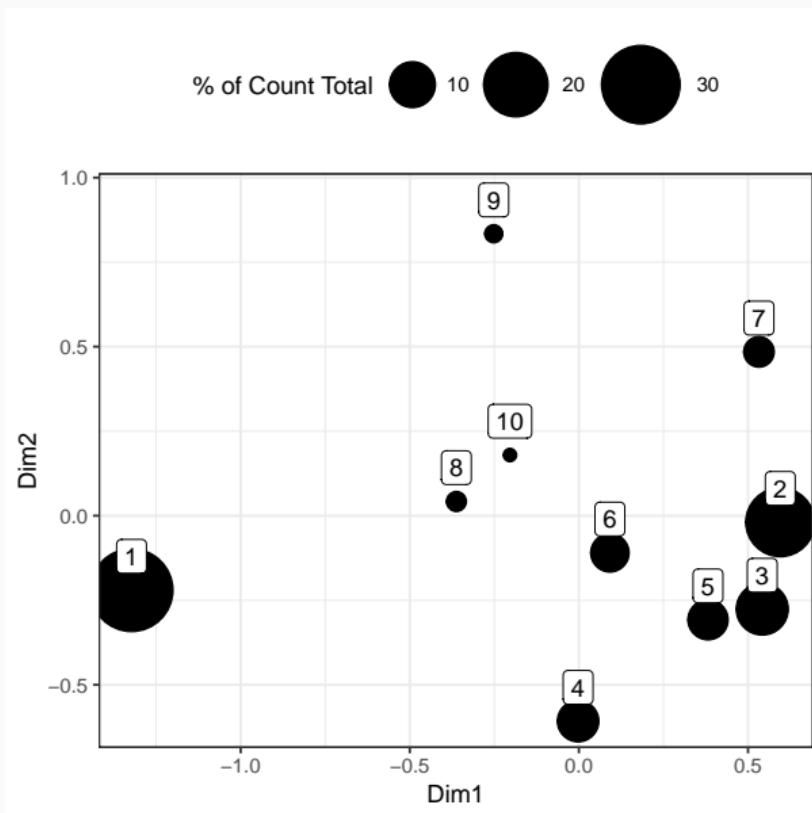
User supplies J – the number of associations *a priori* – $j = \{1, 2, \dots, J\}$

J chosen using AIC, *perplexity*, CV, ...

LATENT DIRICHLET ALLOCATION



LATENT DIRICHLET ALLOCATION — NMDS



LATENT DIRICHLET ALLOCATION — ASSOCIATION 1

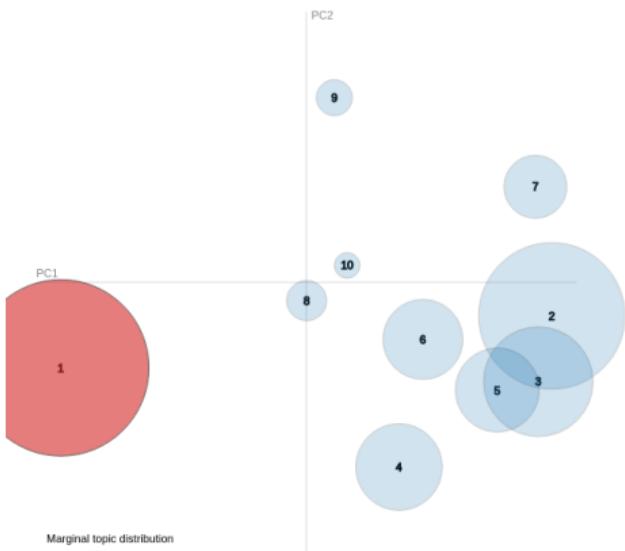
Selected Topic: 0 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:⁽²⁾

$\lambda = 0.5$



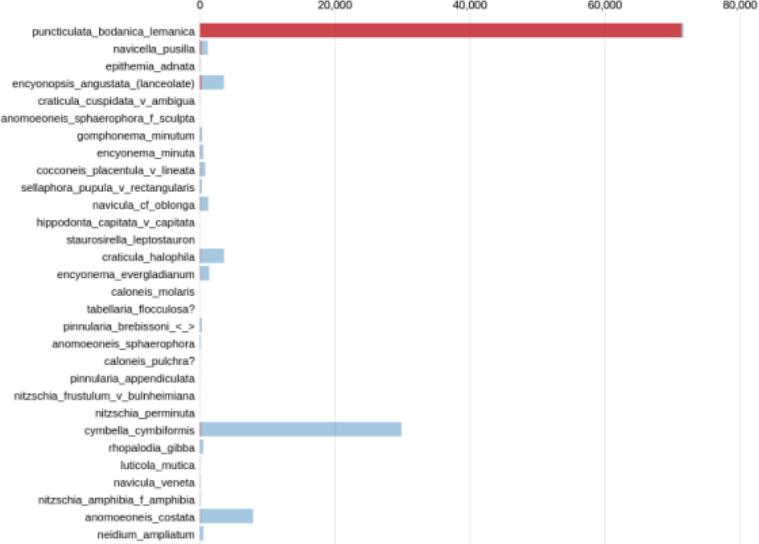
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 1 (33.4% of tokens)



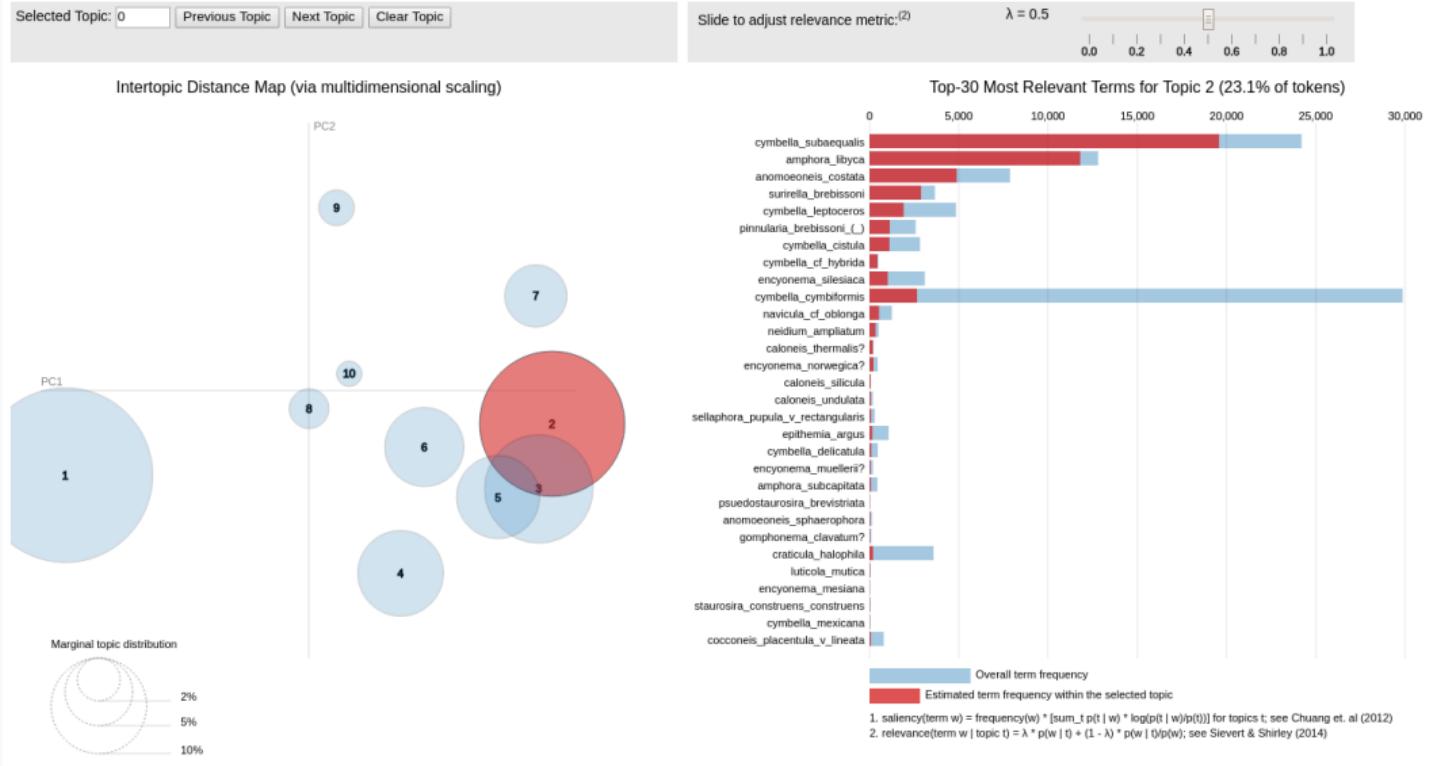
Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et al (2012)

2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t) / p(w)$; see Sievert & Shirley (2014)

LATENT DIRICHLET ALLOCATION — ASSOCIATION 2



LATENT DIRICHLET ALLOCATION — ASSOCIATION 3

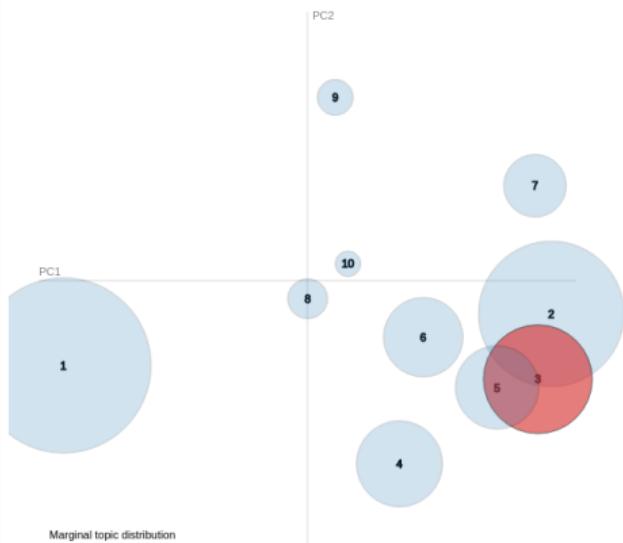
Selected Topic: 0 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:⁽²⁾

$\lambda = 0.5$



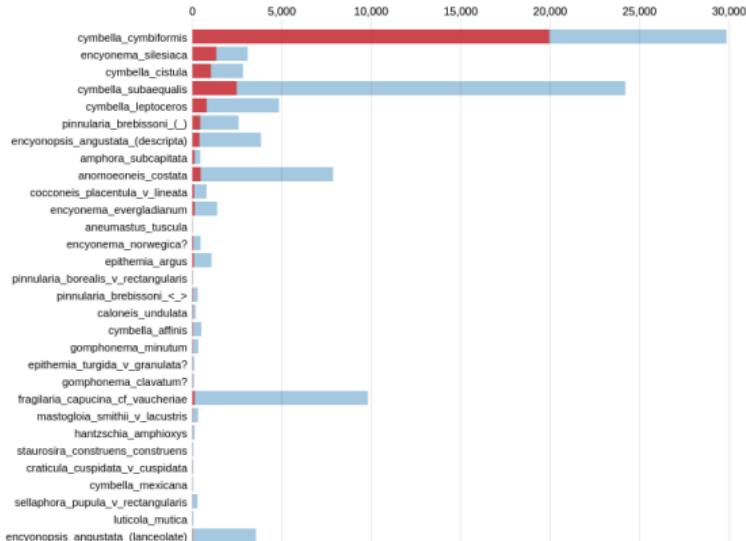
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 3 (12.9% of tokens)



Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = $\text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$ for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)p(w)$; see Sievert & Shirley (2014)

LATENT DIRICHLET ALLOCATION — ASSOCIATION 4

Selected Topic: 0 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:(2)

$\lambda = 0.5$



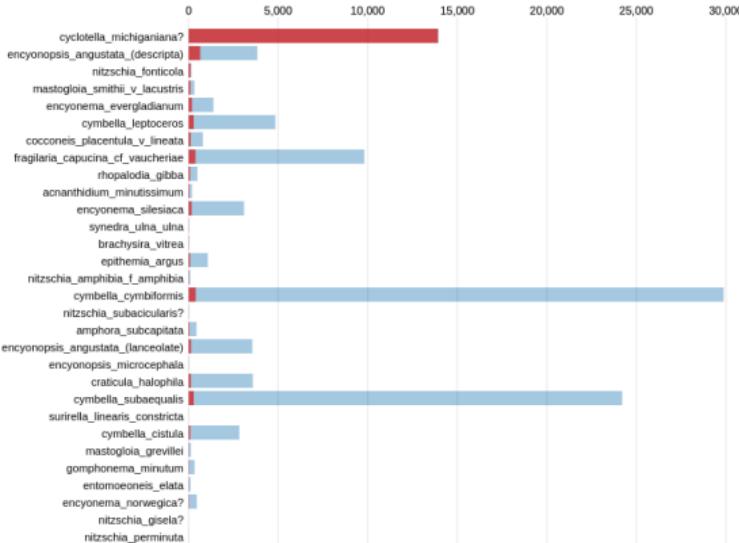
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 4 (8.1% of tokens)



Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)p(t/w)$; see Sievert & Shirley (2014)

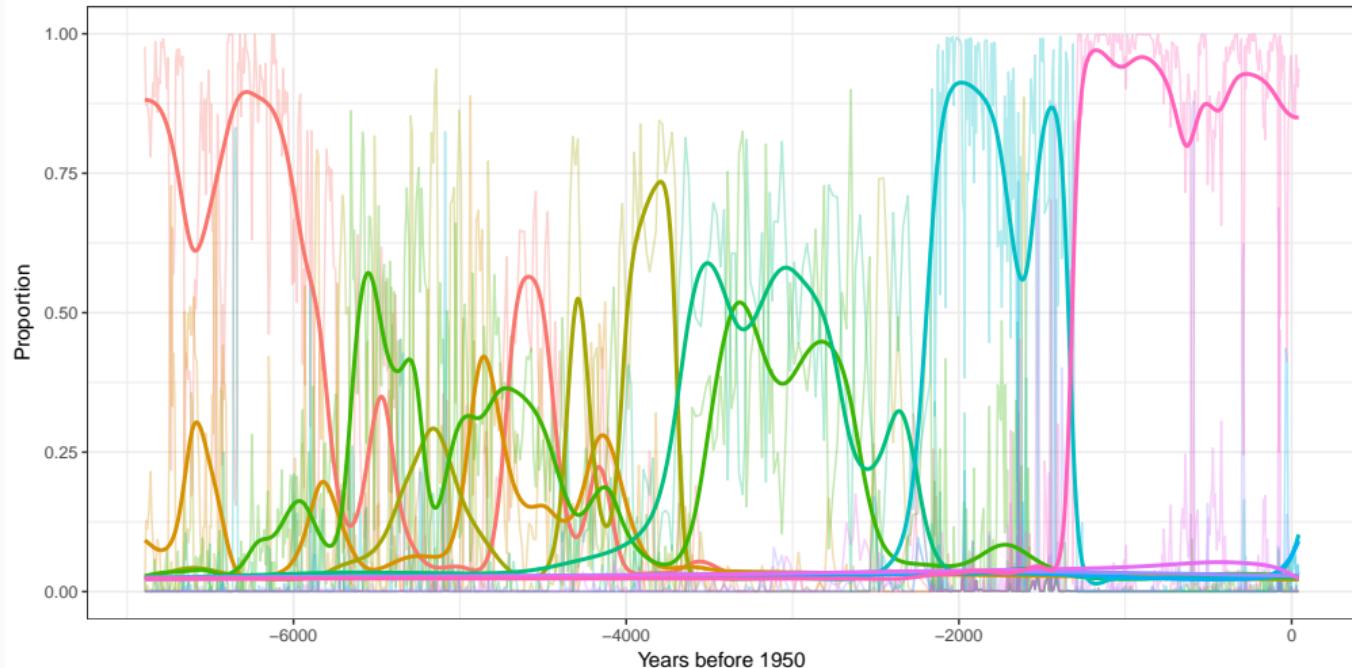
LDA knows nothing about the temporal ordering of the samples

Estimate trends in proportions of species associations using a GAM

- assumes smoothly varying trends
- use adaptive spline to allow for rapid adaptation to changing data
- model each association as $\sim \text{Beta}(\mu, \phi)$

Other methods would also be appropriate: eg. Bayesian Change point model

LATENT DIRICHLET ALLOCATION — TREND ESTIMATION (ADAPTIVE GAM)



CORRELATED TOPIC MODEL

LDA assumes associations of species are **uncorrelated**

Potentially more *parsimonious & realistic* if associations were correlated

2. Proportions of each type in the Skittle community – draw

$$\eta \sim N(\mu, \Sigma)$$

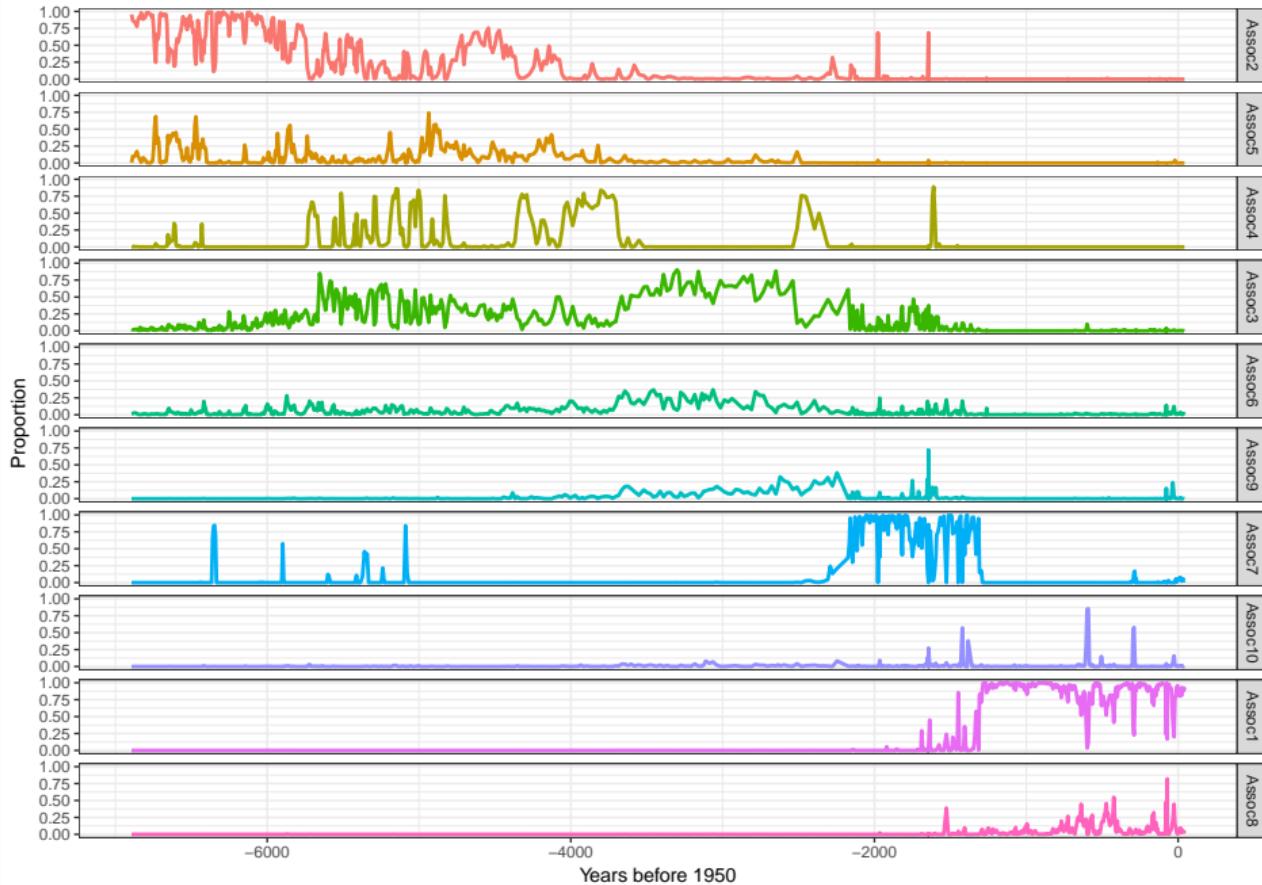
with $\eta \in \mathbb{R}^{J-1}$ and $\Sigma \in \mathbb{R}^{(J-1) \times (J-1)}$

Then transform η_J to proportional scale

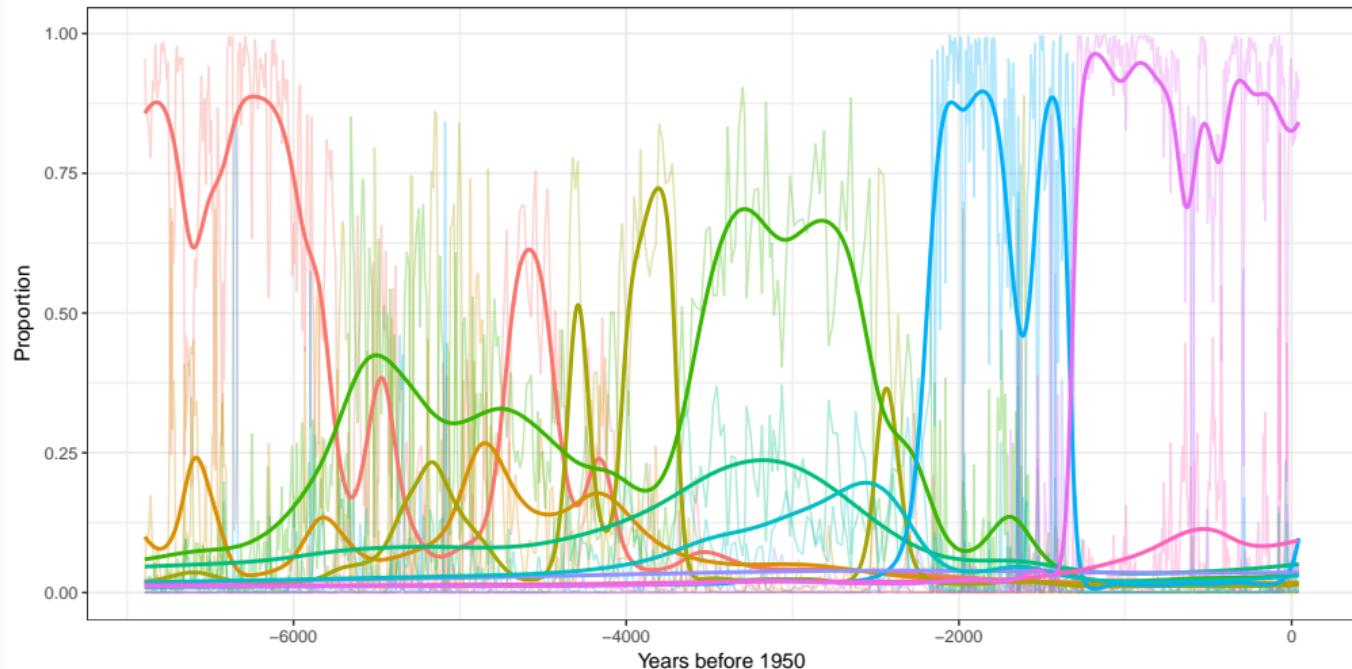
Σ controls the correlation between *associations*

Blei & Lafferty (2007) A correlated topic model of Science. *Ann. Appl. Stat.* 1, 17–35.

CORRELATED TOPIC MODEL



CORRELATED TOPIC MODEL — TREND ESTIMATION (ADAPTIVE GAM)



SUMMARY

LDA & CTM proved well-capable of summarizing the complex community dynamics of Foy Lake

- Reduced 113 taxa to 10 associations of species
- Species associations match closely the expert interpretation of the record
 - make autecological sense also
- The CTM was more parsimonious — removed one rare association
- Estimated trends in proportions of species associations capture
 - smooth, slowly varying trends, and
 - rapid (regime shift?) state change ~ 1.3 ky BP

FUTURE DIRECTIONS

Choosing J is inconvenient

Address this via **Hierarchical Dirichlet Processes** and Bayesian Nonparametrics

- assume J is infinite & put a prior distribution over J

Associations in LDA & CTM are static — distributions are fixed for all samples

- dynamic topic models allow distributions to vary smoothly with time

Many developments in this field:

- *Chinese Restaurant Process*,
- *Indian Buffet Process*, &
- ...

LATENT DIRICHLET ALLOCATION

1. Flavour distribution for j th type of Skittle

$$\beta_j \sim \text{Dirichlet}(\delta)$$

2. Proportions of each type in the Skittle community

$$\theta \sim \text{Dirichlet}(\alpha)$$

3. For each skittle s_i

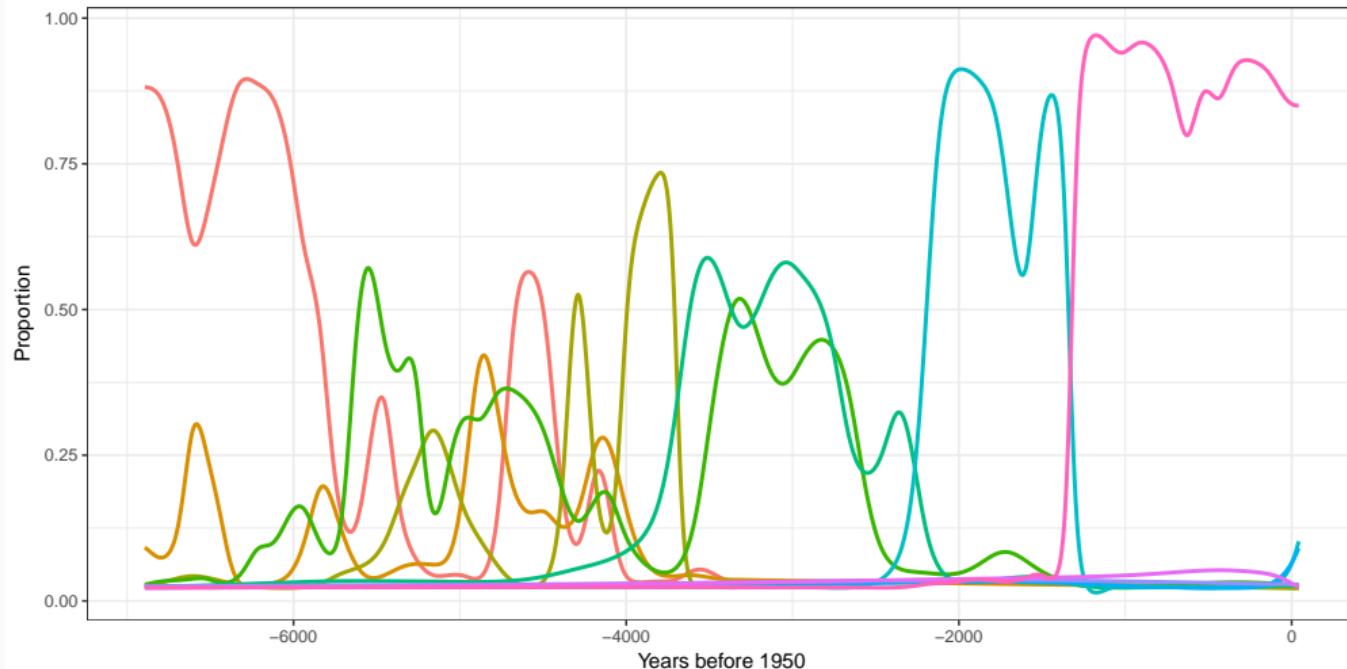
- Choose a pack in proportion

$$z_i \sim \text{Multinomial}(\theta)$$

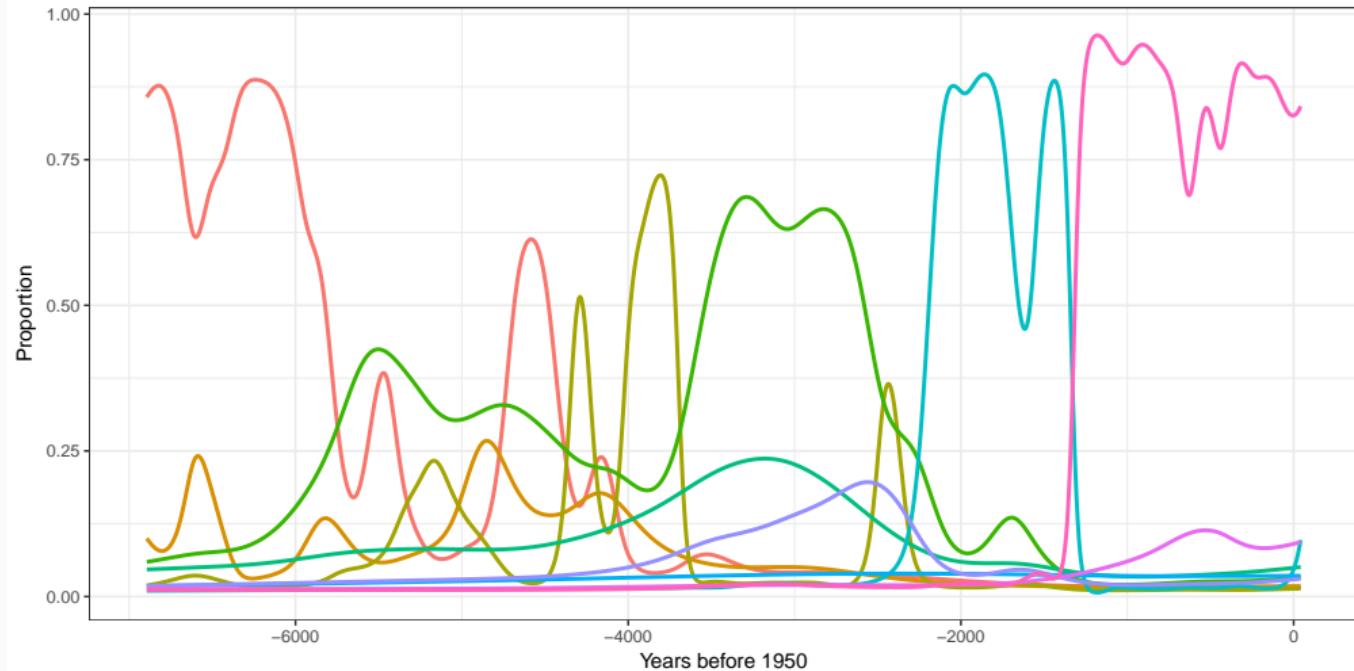
- Choose a flavour from chosen pack with probability

$$p(s_i|z_i, \beta_j) \sim \text{Multinomial}(\delta)$$

LATENT DIRICHLET ALLOCATION — TREND ESTIMATION (ADAPTIVEVGAM)



CORRELATED TOPIC MODEL — TREND ESTIMATION (ADAPTIVE GAM)



INTUITION BEHIND LDA

Latent Dirichlet allocation represents a trade-off between two goals

1. for each sample, allocate its individuals to a **few associations of species**
2. in each association, assign high probability to a **few species**

These are in opposition

- assigning a sample to a single association makes **2 hard** — all its species must have high probability under that one topic
- putting very few species in each association makes **1 hard** — to cover all individuals in a sample must assign sample to many associations

Trading off these two goals therefore results in LDA finding tightly co-occurring species