# Modeling research topics in movement ecology

**Rocío Joo**
ME Boone, S Picardi, VS Romero-Romero, TA Clay, SC Patrick, M Basille
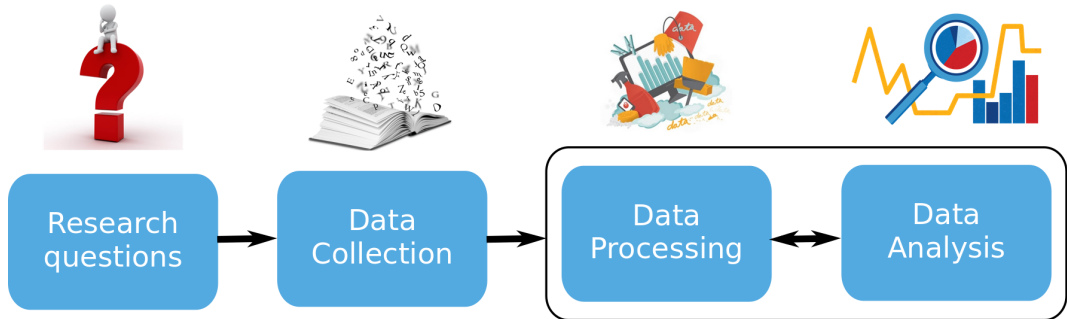
IBC - 2020

✉ rocio.joo@ufl.edu          🐦 rocio_joo

# Text analysis workflow

Research
questions

Data
Collection

Data
Processing

Data
Analysis

— which topics?

— statistical methods?

— software?

.
.
.

in the last decade

Research questions → Data Collection → [ Data Processing ⟷ Data Analysis ]

- which topics?
- statistical methods?
- software?
  .
  .

in the last decade

Research
questions

Data
Collection

Data
Processing

Data
Analysis

( which topics? )  **8007** abstracts

— statistical methods?

— software?

            .
            .
            .

in the last decade

```
Research          Data          Data          Data
questions     →   Collection →  Processing ↔  Analysis
```

— which topics?    **8007** abstracts
— statistical methods?
— software?
        .
        .
        .
in the last decade

— remove
  redundant words

— standard English
  (BrE or AmE)

— lemmatize

— filter out words
  appearing once

Research questions → Data Collection → Data Processing ↔ Data Analysis

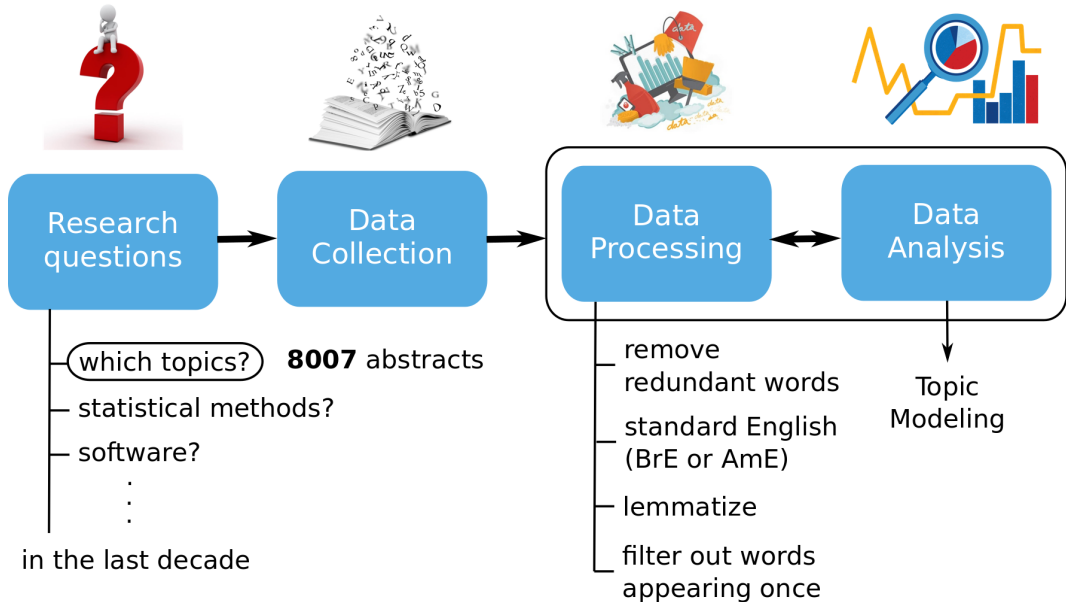- (which topics?)  **8007** abstracts
- statistical methods?
- software?
  .
  .
  .

in the last decade

Data Processing:
- remove redundant words
- standard English (BrE or AmE)
- lemmatize
- filter out words appearing once

Data Analysis → Topic Modeling

Why attend IBC?

## Why attend IBC?

Scientific community

- Community
- Network
- People
- Meeting
- Social

Knowledge

- Learn
- Teach
- Development
- Understand
- Study

## Why attend IBC?

Scientific community

- Community
- Network
- People
- Meeting
- Social

Knowledge

- Learn
- Teach
- Development
- Understand
- Study

*"I can learn new things"*

*"It's nice to see the recent developments in the community"*

*"I am new in the field and want to meet people"*

## Why attend IBC?

Scientific community

- Community
- Network
- People
- Meeting
- Social

Knowledge

- Learn
- Teach
- Development
- Understand
- Study

*"I can learn new things"*

*"It's nice to see the recent developments in the community"*

*"I am new in the field and want to meet people"*

# Topic modeling

Latent Dirichlet Allocation (LDA)

- Bayesian mixture model (Blei et al. 2003; Grün and Hornik 2011)
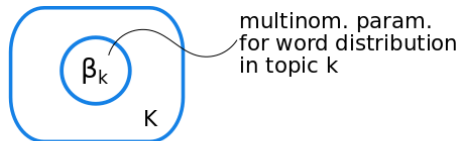
# Topic modeling

Latent Dirichlet Allocation (LDA)

- Bayesian mixture model (Blei et al. 2003; Grün and Hornik 2011)
- From a fixed number $K$ of topics, each topic can be characterized by a multinomial distribution of words with parameter β, drawn from a Dirichlet distribution with param. δ
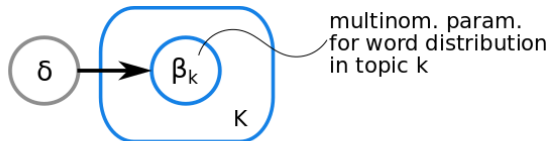
# Topic modeling

Latent Dirichlet Allocation (LDA)

- Bayesian mixture model (Blei et al. 2003; Grün and Hornik 2011)
- From a fixed number $K$ of topics, each topic can be characterized by a multinomial distribution of words with parameter $\beta$, drawn from a Dirichlet distribution with param. $\delta$



multinom. param. for word distribution in topic k

# Topic modeling

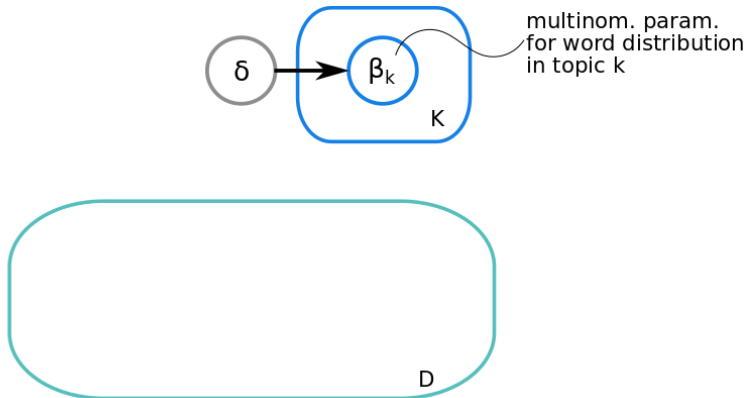Latent Dirichlet Allocation (LDA)

- Bayesian mixture model (Blei et al. 2003; Grün and Hornik 2011)
- From a fixed number $K$ of topics, each topic can be characterized by a multinomial distribution of words with parameter $\beta$, drawn from a Dirichlet distribution with param. $\delta$



multinom. param. for word distribution in topic k
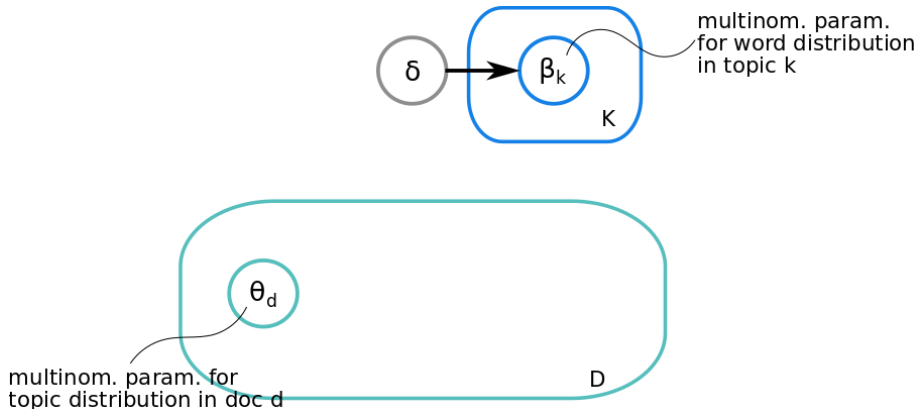
# Topic modeling

Latent Dirichlet Allocation (LDA)

- Each document $\mathbf{d} \in \{1, ..., \mathbf{D}\}$ is composed by a mixture of topics, drawn from a multinomial distribution with parameter $\theta$, which is drawn from a Dirichlet distribution with parameter $\alpha$.



multinom. param. for word distribution in topic k

# Topic modeling

Latent Dirichlet Allocation (LDA)

- Each document $\mathbf{d} \in \{1, ..., \mathbf{D}\}$ is composed by a mixture of topics, drawn from a multinomial distribution with parameter $\theta$, which is drawn from a Dirichlet distribution with parameter $\alpha$.



multinom. param. for word distribution in topic k

multinom. param. for topic distribution in doc d
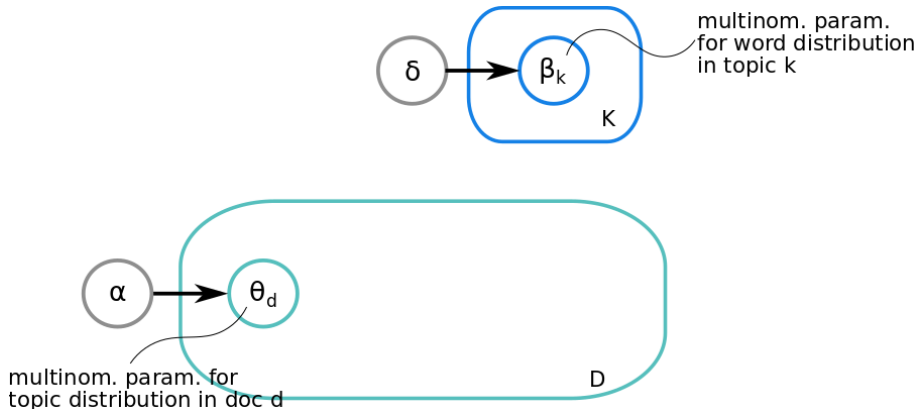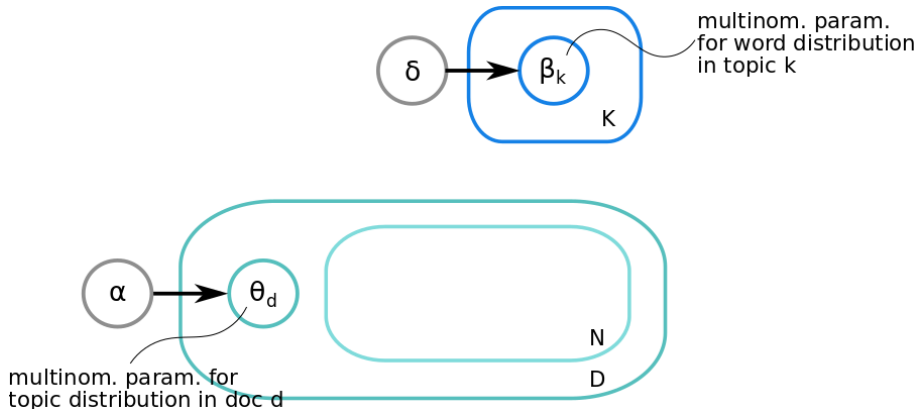
# Topic modeling

Latent Dirichlet Allocation (LDA)

- Each document $\mathbf{d} \in \{1, ..., \mathbf{D}\}$ is composed by a mixture of topics, drawn from a multinomial distribution with parameter $\theta$, which is drawn from a Dirichlet distribution with parameter $\alpha$.



multinom. param. for word distribution in topic k

multinom. param. for topic distribution in doc d

# Topic modeling

Latent Dirichlet Allocation (LDA)

- For each word $\mathbf{w}$ in document $\mathbf{d}$, first a hidden topic $\mathbf{z}$ is selected from the multinomial distribution with parameter $\theta$.



multinom. param. for word distribution in topic k

multinom. param. for topic distribution in doc d

# Topic modeling

Latent Dirichlet Allocation (LDA)

- For each word $\mathbf{w}$ in document $\mathbf{d}$, first a hidden topic $\mathbf{z}$ is selected from the multinomial distribution with parameter $\theta$.
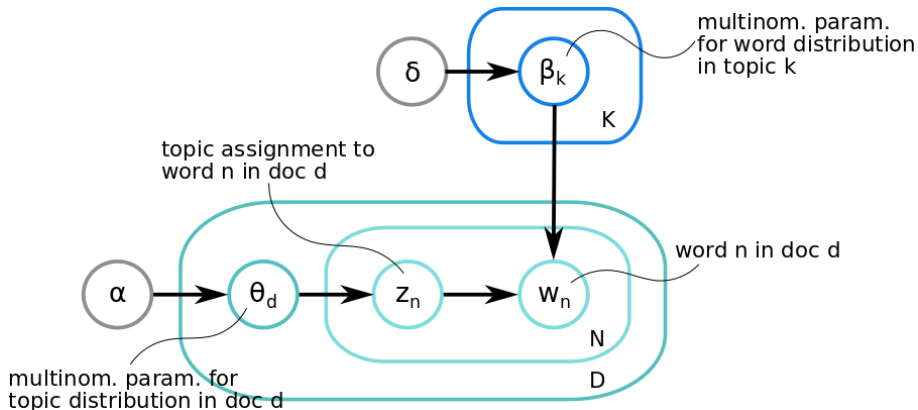
# Topic modeling

Latent Dirichlet Allocation (LDA)

- For each word $\mathbf{w}$ in document $\mathbf{d}$, first a hidden topic $\mathbf{z}$ is selected from the multinomial distribution with parameter $\theta$.
- From the selected topic $\mathbf{z}$, a word is selected based on the multinomial distribution with parameter $\beta$.
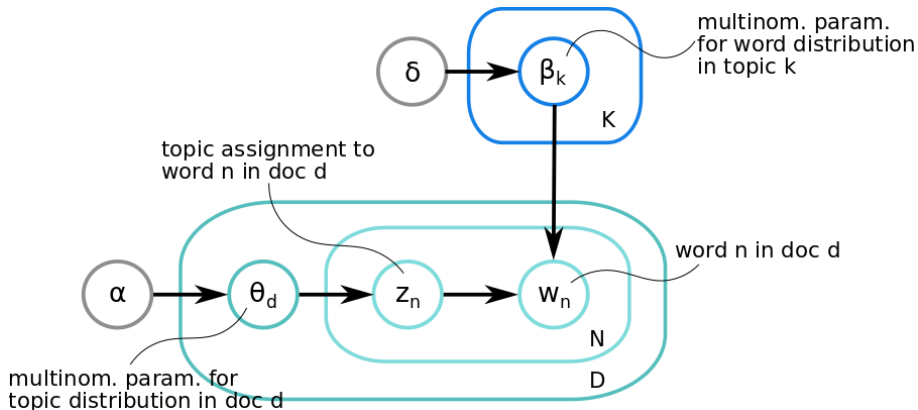
# Topic modeling

Latent Dirichlet Allocation (LDA)

The log-likelihood of a document $\mathbf{d} = \{\mathbf{w_1}, ..., \mathbf{w_N}\}$ is

$$\mathbf{l}(\alpha, \beta) = \log(\mathbf{p}(\mathbf{d} \mid \alpha, \beta)) = \log \int \sum_{\mathbf{z}} \left[ \prod_{\mathbf{n=1}}^{\mathbf{N}} \mathbf{p}(\mathbf{w_i} \mid \mathbf{z_i}, \beta) \mathbf{p}(\mathbf{z_i} \mid \theta) \right] \mathbf{p}(\theta \mid \alpha) \mathbf{d}\theta$$



multinom. param. for word distribution in topic k

topic assignment to word n in doc d

word n in doc d

multinom. param. for topic distribution in doc d

# Topic modeling

Latent Dirichlet Allocation (LDA)

The log-likelihood of a document $\mathbf{d} = \{\mathbf{w_1}, ..., \mathbf{w_N}\}$ is

$$l(\alpha, \beta) = \log(\mathbf{p}(\mathbf{d} \mid \alpha, \beta)) = \log \int \sum_{\mathbf{z}} \left[ \prod_{\mathbf{n=1}}^{\mathbf{N}} \mathbf{p}(\mathbf{w_i} \mid \mathbf{z_i}, \beta) \mathbf{p}(\mathbf{z_i} \mid \theta) \right] \mathbf{p}(\theta \mid \alpha) \mathbf{d}\theta$$

**Inference** on the posterior via Variational Expectation Maximization
(Blei et al. 2003; Grün and Hornik 2011); `topicmodels` R package

**Assumptions:**

- Exchangeability: order of words is negligible
- Topics are uncorrelated
- Number of topics is known (in this study: 15)

Results - movement ecology
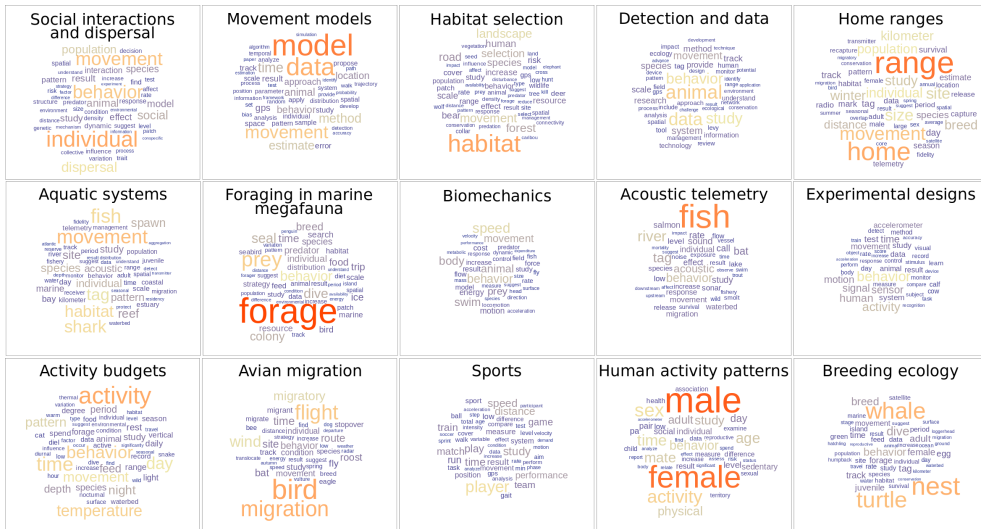
# Topic modeling - LDA

Results - movement ecology

- $\mathbf{E}(\beta \mid \mathbf{z}, \mathbf{w}) \rightarrow$ word distribution per topic $\rightarrow$ label topics

# Topic modeling - LDA

Results - movement ecology

- $\mathbf{E}(\beta \mid \mathbf{z}, \mathbf{w}) \rightarrow$ word distribution per topic $\rightarrow$ label topics

Results - movement ecology

- $\mathbf{E}(\beta \mid \mathbf{z}, \mathbf{w}) \rightarrow$ word distribution per topic $\rightarrow$ label topics
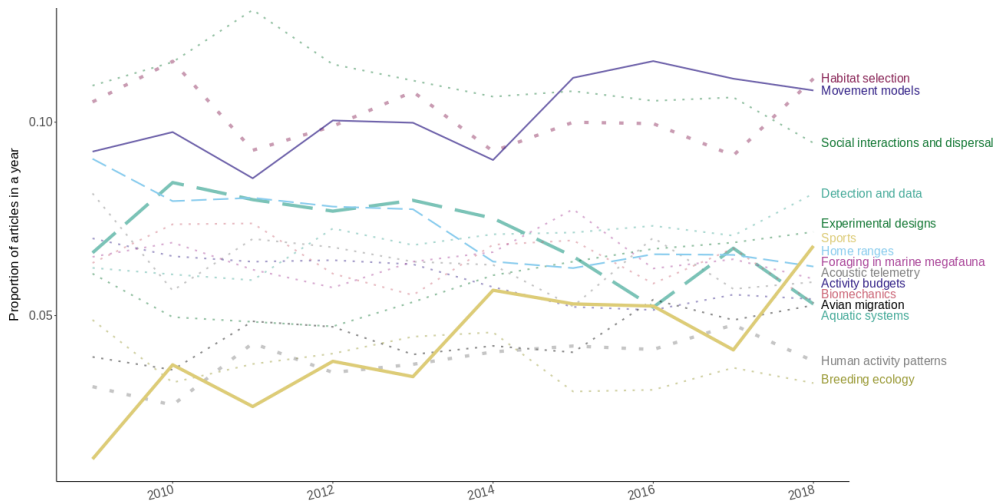
# Topic modeling - LDA

Results - movement ecology

- $\mathbf{E}(\theta_\mathbf{d} \mid \mathbf{z})$: topic distribution per document
- $\sum_\mathbf{d} \mathbf{E}(\theta_\mathbf{d} \mid \mathbf{z_k})$: proxy of prevalence of each topic

# Topic modeling - LDA

Results - movement ecology

- $\mathbf{E}(\theta_\mathbf{d} \mid \mathbf{z})$: topic distribution per document
- $\sum_\mathbf{d} \mathbf{E}(\theta_\mathbf{d} \mid \mathbf{z_k})$: proxy of prevalence of each topic

# Further exploration?

- Topic evaluation: word intrusion
  - Take the highest probability words from a topic
  - Take a high-probability word from another topic and add it
  - Ask humans to identify the word that does not belong

Results shown here: Joo et al. pre-print.

# References

- Blei, David M., Andrew Y. Ng, and Michael I. Jordan. 2003. 'Latent Dirichlet Allocation.' Journal of Machine Learning Research 3 (Jan): 993âĂŞ1022. http://jmlr.csail.mit.edu/papers/v3/blei03a.html.

- Grün, Bettina, and Kurt Hornik. 2011. 'topicmodels: An R Package for Fitting Topic Models.' Journal of Statistical Software 40 (13): 1âĂŞ30. https://doi.org/10.18637/jss.v040.i13.

- Joo, Rocío, Simona Picardi, Matthew E. Boone, Thomas A. Clay, Samantha C. Patrick, Vilma S. Romero-Romero and Mathieu Basille. Pre-print. 'A decade of movement ecology'. arxiv. https://arxiv.org/abs/2006.00110.

Thanks for your attention

✉ rocio.joo@ufl.edu  🐦 rocio_joo