# Unveiling hidden semantic structures of corpora using Topic Models

Simón Roca and Jerónimo Arenas

May 14, 2018

Universidad Carlos III de Madrid, Dpto. TSC

Presented by Simón Roca (sroca@ing.uc3m.es) in Machine Learning Group, UC3M.

## Outline

## Outline

# Structured corpus vs. Unstructured corpus

Some examples: project forms, job offers, course guides, patents, articles...

- Inputs: BoW corpus, hyperparameters...
- Outputs: document-topic proportion matrix, topic-word proportion vectors...

## Our proposal

- Split the documents into paragraphs.
- New plate/variable: let the paragraphs be semantic or background.
- Use different LDA parameters for each kind.

## Outline

Modeling General and Specific Aspects of Documents with a Probabilistic Topic Model

- Three sets of topics: classic topics, corpus specific topic, and document specific topic.

(b) HIERSUM Graphical Model

Exploring Content Models for Multi-Document Summarization

- Used for summarizing documents.
- They assume constant topic across a single sentence.

Correlated Topic Models

- Logistic Normal distribution as a prior: correlation among topics can be obtained.
- Non Conjugacy, but permits clustering topics.

## Previous approaches: Hierarchies among topics



(b) Four-Level PAM

Pachinko Allocation: DAG-Structured Mixture Models of Topic Correlations

- It allows learning topics, sub-topics, super-topics...
- Due to its tree structure (DAG), it permits obtaining correlations.

Hidden Topic Markov Models

- Each word in a sentence comes from the same topic.
- Transition between sentences modeled with a Markov Chain.

Figure 1: The *senLDA* model. The words $w$ of a sentence share the same topic $z$.

On a topic model for sentences

- Each word in a sentence comes from the same topic.
- Modeled as an extra plate in model.

## Is our model different?

- We propose a model combining some of the previous approaches.
- However, that's not the main contribution.
- Why not selecting in which parts useful topics may be learned, and then looking there for best quality topics?

## Outline

Figure 2: Paragraph LDA graphical model

# Generative Model

1. $\phi_1 \sim Dir_{V_1}(\beta_1)$

2. $\phi_0 \sim Dir_{V_0}(\beta_0)$

3. For each document,

    (a) $\theta_1 \sim Dir_{K_1}(\alpha_1 << 1)$

    (b) $\theta_0 \sim Dir_{K_0}(\alpha_0 \geq 1)$

    (c) For each paragraph,

- $\psi_x \sim Dir_2(\gamma)$
- $x \sim Ber(\psi_x)$
  - if $x = 1$, for each word:
    * $x_w \sim Ber(m)$
    * if $x_w = 1$:
      · $z \sim Mult(\theta_1)$
      · $w \sim Mult(\phi_{1,z})$
    * if $x_w = 0$:
      · $z \sim Mult(\theta_0)$
      · $w \sim Mult(\phi_{0,z})$
  - if $x = 0$, for each word:
    * $z \sim Mult(\theta_0)$
    * $w \sim Mult(\phi_{0,z})$

## Inference

$$p(\vec{z}, \vec{x}, \vec{w}) = p(\vec{w}|\vec{z}, \vec{\beta})p(\vec{z}|\vec{x}, \vec{\alpha})p(\vec{x}|\gamma) \qquad (1)$$

$$p(\vec{x}|\gamma) = \frac{\Delta(\vec{n}_d + \gamma)}{\Delta(\gamma)} \qquad (2)$$

$$p(z_i = k|\vec{z_{\neg i}}, \vec{w}) \propto (n_{m,\neg i}^{(t)} + \alpha_k)\frac{n_{k,\neg i}^{(t)} + \beta_t}{\sum_{t=1}^{V} n_{k,\neg i}^{(t)} + \beta_t} \qquad (3)$$

$$p(x_p = s|\vec{w}, \vec{x}_{\neg p}, \vec{z}, \alpha, \beta, \gamma) \propto \frac{\prod_{w \in p}(n_{s,\neg p}^{(w)} + \gamma)...(n_{s,\neg p}^{(w)} + \gamma + (n_{s,p}^{(w)} - 1))}{(\sum_{w \in V}(n_{s,\neg p}^{(w)} + \gamma))...(\sum_{w \in V} n_{s,\neg p}^{(w)} + \gamma + (n_{s,p}^{(w)} - 1))} \qquad (4)$$

## Outline

## Evaluation metrics: Histogram intersection distances

For matching real topics and predicted ones in synthetic dataset.
Checking how much they differ.

$$exp(-\alpha * \sum n_{j=1} min(I_j, M_j)) \tag{5}$$

## Evaluation metrics: Topic coherence

- Perplexity (held-out probability) is not the best choice for topic quality/coherence.
- Some measurements based on PMI (word co-ocurrences), WordEmbeddings (word vectors), high correlation with human judgement.
- Learned on reference corpus, strange behaviour with specific vocabulary.
- David Mimno stands that PMI can be obtained in the same corpus.

**Experiments on synthetic dataset: some numbers**

When the generative model is met, is our method worth?

| Attribute | Docs(test) | Paragraphs | Total Words | $K_0(K_1)$ | $V$ | $\alpha_0(\alpha_1)$ | $\beta_0(\beta_1)$ |
|---|---|---|---|---|---|---|---|
| Value | 3000(500) | ? | ? | 10 (30) | 5000 | 2 (0.1) | 0.1(0.1) |

**Table 1:** Synthetic dataset generated using our generative model

Noisy: variable proportion of background words in semantic paragraphs.

It seems so!

## Experiments on real datasets

After several runs in Patents, SCIELO and NIPS:

- Good intuition about learned topics...
- ...But referenced corpus coherence doesn't claim so.
- In some corpus there are too many semantic paragraphs.

## Experiments on real datasets: Patents

Some background topics:
['temperature', 'gas', 'liquid', 'water', 'heat', 'heating', 'pressure', 'oil', 'fluid', 'chamber']
['compound', 'reaction', 'mixture', 'mmol', 'substituted', 'solution', 'atom', 'formula', 'acid', 'stirred']
['page', 'subject', 'action', 'sequence', 'active', 'report', 'activity', 'factor', 'total', 'property']

Some semantic topics: ['channel', 'transmission', 'antenna', 'receiver', 'symbol', 'transmit', 'transmitted', 'transmitter', 'resource', 'carrier']
['server', 'client', 'call', 'network', 'web', 'request', 'status', 'telephone', 'manager', 'local']
['frequency', 'phase', 'filter', 'digital', 'noise', 'pulse', 'band', 'amplitude', 'clock', 'gain']

Proper experiments:

- Feasible: running parLDA in these datasets and compare some topics to those obtained from classic LDA.
- Just in time: small validation process on hyperparameters based on topic coherence.
- ...? Deadline on Friday, 8pm.

## Outline

## Conclusions

- Our inference performs well enough when the generative model is met.

- In that scenario, identifies better semantic paragraph than other methods, even without labels.

- On real datasets, it finds different paragraphs, learning reasonable topics.

## Future Work

- Deep analysis of the quality of these new topics.

- Analysis of topic coherence measurements.

- Getting richer models.

- ...