# Automatic Labeling for Multinomial Topic Models

Problem: meaningful labels for discovered topics.

w/o being subjective through manual labeling *

Solution: optimization problem involving minimizing

Kullback-Leiber divergence while maximizing mutual

information between a **label** and a topic **model**.

Href:

$\longrightarrow$ Try to minimize distance b/w label "$l$" and topic "$\theta$"

Kullback-Leiber Definition:

$$\text{Divergence} \leadsto D_{KL}(P\|Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right) \neq D_{KL}(Q\|P)$$

$\{trump, ukraine, zelensky, putin\} = \theta$

then $D_{KL}(\theta\|l_1) > D_{KL}(\theta\|l_2)$

$L = \begin{cases} l_1 = \text{"war on ukraine"} \\ l_2 = \text{"sports game"} \\ l_3 = \text{"cooking show"} \end{cases}$

$\downarrow$ Probability Space

Href:

Context: topic model $\theta$ in text collection $C$ is a probability

distribution of words $\omega$ in vocab set $V$:

$$\sum_{\omega \in V} P(\omega | \theta) = 1$$

topic label : sequence of words which is semantically meaningful and covers the latent meaning of topic $\theta$.

Relevance Score : $S(\ell, \theta)$ measures Semantic $\overset{\text{relative to long or}}{\underset{\text{topic}}{\nearrow}}$ similarity between the label and the topic model.

To generate labels that are understandable, semantically relevant, discriminative across topics, and of high coverage of each topic, we first extract a set of candidate labels in a preprocessing step.*

$\hookrightarrow$ Best labels need to be phrases present within the collection $C$

$\underset{\nearrow}{\quad}$ "Candidate label"

Phrase Generation:

$\hookrightarrow$ Chunking / Shallow Parsing $\sim$ Aims to identify short phrases in text using parts of speech tags to make decisions of chunking according to some grammar.

• look for chunks / phrases frequently appearing in text as candidate labels

↪ Ngram Testing : extract meaningful phrases from word

ngrams based on Statistical tests.

    ↪ if words in an ngram co-occur w/ eachother,

    the n-gram is more likely to be an n-word phrase.

Href : https://medium.com/analytics-vidhya/generating-meaningful-phrases-from-unstructured-news-data-d4e217a7da43

## Semantic Relevance Scoring

↪ The zero-order Relevance

any reasonable measure of relevance b/w $\theta$ and $\ell$

should compare $\ell$ to the distribution of words that

define $\theta$.

$\theta$ is a    $Score = \log \dfrac{P(\ell|\theta)}{P(\ell)} = \sum\limits_{0 \le i \le m} \log \dfrac{P(\omega_i|\theta)}{P(\omega_i)}$

topic

     ↪ $\omega_i$ is a

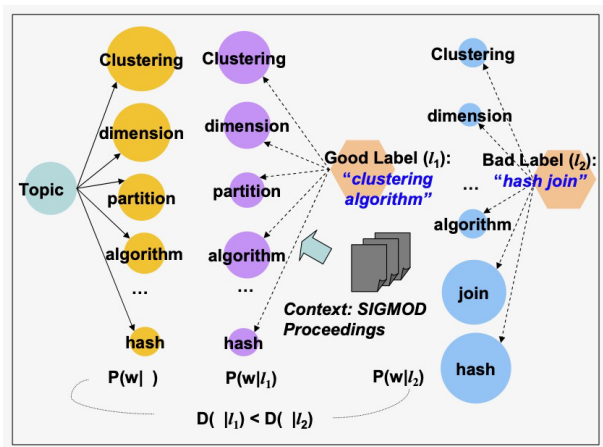    $\ell$ is a phrase      where $\ell = \omega_0, \omega_1 \dots \omega_m$    word

• a phrase containing more "important" (high $P(\omega|\theta)$)

words in the topic distribution is assumed to be

a good label.

• Does not look at context information from reference collection

↳ The first-order Relevance

• The semantics of a topic model should be interpreted in a context.

• A natural context to interpret a topic is the original collection from which the topic model is extracted

◦ Represent a candidate label also with a multinomial distribution of words. → $\{P(w|l)\}$ decided by $l$

can measure closeness of $\{P(w|l)\}$ and $\{P(w|\theta)\}$ using

Kullback-Leiber divergence

$$D(\theta \| l)$$



• To use this relevance score, we need to approximate distribution $\{P(w|l)\}$ ↝ include a context collection $C$.
↳ substitute $\{P(w|l)\}$ w/ $\{P(w|l,C)\}$

$\hookrightarrow$ Relevance Scoring Function

$$\text{Score}(l, \theta) = -D(\theta \| l) = -\sum_{\omega} p(\omega | \theta) \log \frac{p(\omega | \theta)}{p(\omega | l)}$$

$$= -\sum_{\omega} p(\omega | \theta) \log \frac{p(\omega | c)}{p(\omega | l, c)} - \sum_{\omega} p(\omega | \theta) \log \frac{p(\omega | \theta)}{p(\omega | c)} - \sum_{\omega} p(\omega | \theta) \log \frac{p(\omega | l, c)}{p(\omega | l)}$$

$$= \sum_{\omega} p(\omega | \theta) \log \frac{p(\omega, l | c)}{p(\omega | c) p(l | c)} - D(\theta \| c) - \sum_{\omega} p(\omega | \theta) \log \frac{p(\omega | l, c)}{p(\omega | l)}$$

$$= \sum_{\omega} p(\omega | \theta) \, PMI(\omega, l | c) - D(\theta \| c) + Bias(l, c)$$

$\downarrow$        $\searrow$ Bias of using

Divergence     Context $C$ to infer

• Can ignore     b/w the topic     the semantic relevance

bias if all candidate    and the labeling    of $l$ and $\theta$.

labels are generated    context.

from collection $C$.     • identical if      • can be utilized

all candidate labels     to incorporate priors

came from same     of candidate labels.

collection

• First component can be written as the expectation of pointwise

mutual information b/w $l$ and the terms in the topic model

given the context $C \, E_\theta (PMI(\omega, l | c)))$.

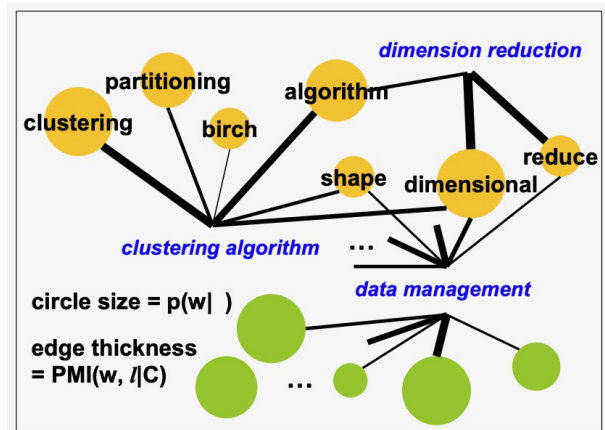• PMI of $x, y$ is correlation b/w events $x$ and $y$

$$= \log \left( \frac{p(x, y)}{p(x) p(y)} \right)$$

- PMI $(w, l \mid C)$ can all be precomputed independently of the topic models to be labeled.

\# Use Leplace smoothing to filer out "$P(w, l \mid C) = 0$"-words in order to avoid PMI $(w, l \mid C)$ being undefined

⤷ Intuitive interpretation

- can construct weighted graph where each node is either a term in a topic model or a candidate label.



- Each edge blw label and topical term measured using PMI$(w, l \mid C)$
- The weight of each node indicates the importance of the term to this topic. We. The eweight of each edge indicates how strongly the label and the term are semantically associated.

- The scoring function $E_\theta (PMI(w, l \mid C))$ would rank a label node higher if it generally has stronger semantic relation to these important topical words.

High Coverage Labels