

# Precision-Recall Balanced Topic Modelling

Seppo Virtanen<sup>1</sup> and Mark Girolami<sup>1,2</sup>

University of Cambridge<sup>1</sup>

The Alan Turing Institute<sup>2</sup>

## Summary

Topic models are becoming increasingly relevant probabilistic models for dimensionality reduction of text data, inferring topics that capture meaningful themes of frequently co-occurring terms. In this work, we present:

- New insights into topic modelling from an information retrieval perspective.
- New evaluation measures for topic modelling based on the precision-recall trade-off.
- Novel topic model combined with an efficient inference algorithm that allows the user to balance between contributions of precision and recall, inferring more coherent and meaningful topics.

## Topic models are recall-biased

$M$  documents  $\mathbf{y}_m$ , where  $m = 1, \dots, M$ , such that  $y_{m,d}$ , where  $d = 1, \dots, D$ , denotes a frequency of the  $d$ th term in the vocabulary for the  $m$ th document.  $N_m$  individual words for the  $m$ th document are denoted as  $w_{m,n} \in \{1, D\}$ , where  $n = 1, \dots, N_m$ .

Topic models assume multinomial likelihood

$$\mathcal{L}_m = \prod_{d=1}^D q_{m,d}^{y_{m,d}}$$

where  $\mathbf{q}_m \in \Delta^D$  denotes an unknown expectation parameter of the multinomial distribution, satisfying  $q_{m,d} \geq 0$  and  $\sum_{d=1}^D q_{m,d} = 1$ . The expectations  $\mathbf{q}_m$  decompose as a linear decomposition

$$\mathbf{q}_m = \sum_{k=1}^K \boldsymbol{\eta}_k \theta_{m,k},$$

where  $\boldsymbol{\eta}_k \in \Delta^D$ , for  $k = 1, \dots, K$ , correspond to **topics** and  $\theta_m \in \Delta^K$  to **topic proportions**. Further, assume

$$\boldsymbol{\eta}_k \sim \text{Dirichlet}(\gamma \mathbf{1}), \quad \boldsymbol{\theta}_m \sim \text{Dirichlet}(\boldsymbol{\alpha}).$$

Define an empirical word occurrence distribution

$$p_{m,d} = y_{m,d}/N_m.$$

Inferring  $\mathbf{q}_m$  that is closer to  $\mathbf{p}_m$  leads to more accurate topics. The mean multinomial log likelihood relates to the KL-divergence between empirical and latent word distributions,

$$KL(\mathbf{p}_m, \mathbf{q}_m) = \sum_d (p_{m,d} \log p_{m,d} - p_{m,d} \log q_{m,d}).$$

The divergence is sensitive to the contribution of *misses*, corresponding to terms for which  $\mathbf{p}_m$  are *large* but the corresponding  $\mathbf{q}_m$  are *small*, and, thus, closely relates to the concept of *recall*. Even though, these topics emphasise recall, they may have very low *precision*, containing intruder terms that capture false similarities.

## Information retrieval aspect

Based on the retrieval model  $\mathbf{q}_m$ , the goal is to retrieve co-occurring terms. Here, the  $\mathbf{p}_m$  represent relevances (that is, empirical co-occurrences) and  $\mathbf{q}_m$  should be similar to the  $\mathbf{p}_m$ , avoiding errors. Define two classes of errors, misses and false positives: Terms for which

- 1  $\mathbf{p}_m$  are large but  $\mathbf{q}_m$  are small correspond to **misses**, and
- 2  $\mathbf{q}_m$  are large but  $\mathbf{p}_m$  are small correspond to **false positives**.

Concepts of recall and precision may be quantified with the directed KL divergences, because  $KL(\mathbf{p}_m, \mathbf{q}_m)$  emphasises misses and the reversed divergence  $KL(\mathbf{q}_m, \mathbf{p}_m)$  emphasises false positives.

Consider maximum entropy distributions for  $\mathbf{p}$  and  $\mathbf{q}$  that take uniform values over the support of the distributions, denoted as  $P$  and  $Q$ , respectively, whereas the remaining values are  $\epsilon$ -close to zero ( $\epsilon \approx 0$ ).

$$KL(\mathbf{p}^*, \mathbf{q}^*) = C + \frac{|P \cap Q|}{|P|} \log \epsilon$$

is proportional to *standard* recall, proportion of relevant terms that are retrieved.

$$KL(\mathbf{q}^*, \mathbf{p}^*) = C + \frac{|P \cap Q|}{|Q|} \log \epsilon$$

is proportional to *standard* precision, proportion of retrieved terms that are relevant.

Because of the connections, we may interpret the directed divergences as generalisations of the concepts of recall and precision for continuously-valued grades of relevances.

## Model

Introduce  $0 \leq \lambda \leq 1$  that trade-offs the contributions of recall and precision, capturing term co-occurrences and avoiding false similarities. Assume term-specific assignment variables,

$$x_{m,n} \sim \text{Bernoulli}(\lambda), \quad c_{m,n} \sim \text{Categorical}(\boldsymbol{\theta}_m).$$

If  $x_{m,n} = 0$ , with probability  $1 - \lambda$ ,

$$w_{m,n} \sim \text{Categorical}(\boldsymbol{\eta}_{c_{m,n}}),$$

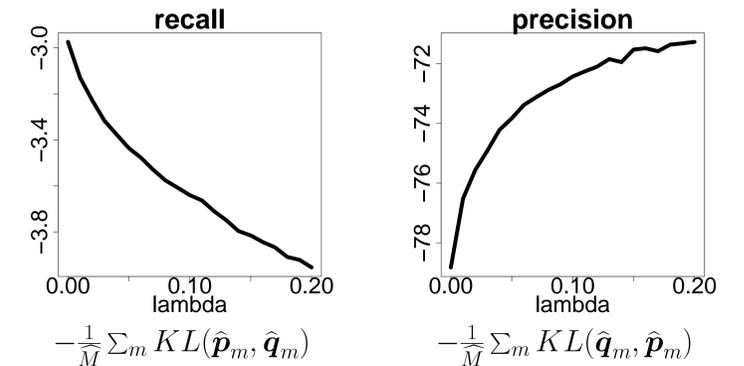
otherwise the term is explained by the  $\mathbf{p}_m$ , with probability  $\lambda$ .

Collapsed Gibbs sampling for inference: For  $w_{m,n} = d$ ,

$$p(c_{m,n} = k, x_{m,n} = 0) \propto \frac{N_{k,m}^{-(w_{m,n})} + \alpha_k}{\sum_{k'} N_{k',m}^{-(w_{m,n})} + \sum_{k'} \alpha_{k'}} \times \frac{G_{k,d}^{-(w_{m,n})} + \gamma}{\sum_{d'} G_{k,d'}^{-(w_{m,n})} + \gamma D}$$

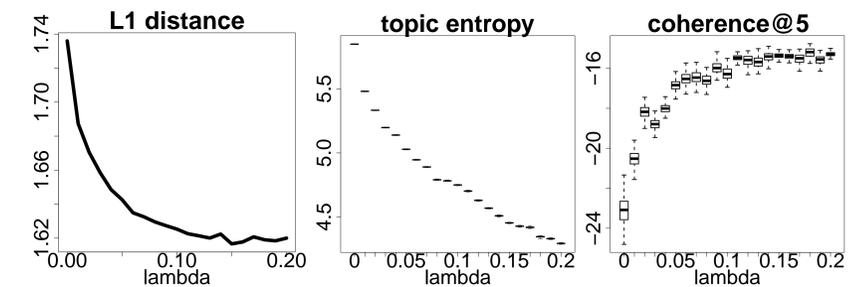
$$p(x_{m,n} = 1) \propto \frac{\lambda}{1 - \lambda} p_{m,d}.$$

## Results



Precision-recall trade-off:

- Recall is maximal for  $\lambda = 0$  and decreases for increasing  $\lambda$ .
- Precision increases for increasing values for  $\lambda$ .



Precision is positively associated with the other performance measures, excluding recall. For increasing  $\lambda$ , and higher precision,

- mean  $\ell_1$  distances between  $\hat{\mathbf{p}}$  and  $\hat{\mathbf{q}}$  become smaller,
- the topics become more sparse, and
- the semantic coherences for the topics increase.

## Acknowledgements

The authors were supported by the EPSRC grant EP/P020720/1, Inference COmputation and Numerics for Insights into Cities (ICONIC), <https://iconicmath.org/>

