

Relevant work

Summaries

[1] SECTOR: A Neural Model for Coherent Topic Segmentation and Classification

Introducing SECTOR – to support machine reading systems – a model that **segments documents into coherent sections and assigns topic labels** to them. In addition, contributes [WikiSection](#), a publicly available data set with 242k labelled sections in English and German from two distinct domains: diseases and cities.

The tasks of topic modelling, text segmentation and text classification are often done separately; SECTOR unifies them. Topic modelling commonly is about extracting the semantic content of entire documents, not segments – mentioned are LDA models ([towardsdatascience example](#)), neural topic embeddings, RNNs. Text segmentation has been done by a wide variety of methods, most recently: LDA, word embeddings, LSTM networks, CNNs. Text classification detects topics on text chunks (paragraphs, sentences) using machine learning methods, such as: SVMs, shallow/deep NNs, CNNs.

- Input →
 - Document $\mathbf{D} = \langle \mathbf{S}, \mathbf{T} \rangle$
 - N sentences $\mathbf{S} = [s_1, \dots, s_N]$
 - Empty segmentation $\mathbf{T} = \emptyset$ as input.
- Task →
 - Split \mathbf{D} into sections $\mathbf{T} = [T_1, \dots, T_M]$
 - $T_j = \langle S_j, y_j \rangle$
 - S_j is a sequence of coherent sentences $S_j \subseteq S$
 - y_j describes the common topic for each sentence in T_j

SECTOR architecture:

1. **Sentence encoding:** transform each plain text sentence into fixed-size sentence vector \mathbf{x}_k , serving as input into the neural network layers. 3 methods:
 - Bag-of-Words Encoding as baseline.
 - [Bloom Filter](#) Embedding, which compresses things.
 - Sentence Embeddings, distributional representation based on pre-trained word2vec embeddings.
2. **Topic embedding:** using 2 layers of LSTM (bidirectional, with forget gates) to produce a “dense distributional representation of latent topics for each sentence”.
3. **Topic classification:** class labels decoded using one-hot or bag-of-words (for single or multi-labelled sections respectively).
4. **Topic segmentation:** leverage the information encoded in the topic embedding and output layers to segment the document and classify each section.
 - *Newline* embedding as a baseline: split sections at every sentence break and then merge sections that at least share one label in the top-2 predictions.
 - Using deviation of topic embeddings: calculate magnitude of *embedding deviation* (emd) per sentence, peaks of emd are used as starts of a new section.
 - Addition to the previous method: *bidirectional embedding deviation* (bemd).

Results were good; bemd best performance; bloom filters on par with sentence embeddings but is slower.

[2] Text Segmentation as a Supervised Learning Task

Previous work on text segmentation focused on unsupervised methods such as clustering or graph search, due to the lack of quality (real-world & varied) labelled data. In this work, text segmentation is formulated as a supervised learning problem, and presented is a large new dataset, *Wiki-727K*, for text segmentation that is automatically extracted and labelled from Wikipedia.

- Document is \mathbf{x} , represented as:
 - sequence of N sentences s_1, \dots, s_n .
 - label $y = (y_1, \dots, y_{n-1})$, represents the segmentation of the document by $n-1$ binary values, where y_i denotes whether s_i ends a segment.
- Model exists of 2 sub-networks:
 - Lower-level is two-layer bidirectional LSTM, which takes a sentence as input and generates its sentence representation/embedding.
 - Higher-level is two-layer bidirectional LSTM, which takes a sequence of sentence embeddings and tries to predict the segmentation probability for each embedding.

Fairly good results, but the model's performance also depends on the kind of text it tries to predict on (e.g., texts from the chemistry domain are more confusing if it hasn't trained enough with similar texts).

[3] Topic segmentation in ASR transcripts using bidirectional RNNs for change detection

Approaches topic segmentation (in speech recognition transcripts) by **measuring lexical cohesion using bidirectional RNNs (LSTM)** – lexical cohesion refers to analysing the lexical (=semantic) distribution across a document and denoting boundaries in areas of low cohesion.

1. (Pre-trained) word embeddings are put into both the forwards and backwards LSTM-RNNs.
2. The *hidden layer activations* from these two networks are transformed using a fully connected feed forward network.
 - a. Forwards output compares similarity between topic context *until* the current word.
 - b. Backwards output compares similarity between topic context *following* the current word.
3. Finally, a dot product of the transformed forwards and backwards outputs are put through an output function. Output function is either:
 - a. Softmin + Cross-entropy cost for training or
 - b. Flipped sigmoid + binary cross-entropy cost for training.
4. The result represents the topic context similarity across all words.

Training & validation datasets are obtained by concatenating 2 or more randomly picked news articles; objective becomes to mark the boundary between articles. Testing was done on both concatenated videos and longer programs containing multiple segments.

[4] A More Effective Sentence-Wise Text Segmentation Approach Using BERT

The proposed system utilizes **Bidirectional Encoder Representations from Transformers (BERT)** as an encoding mechanism, which feeds to several downstream layers with a final classification output layer.

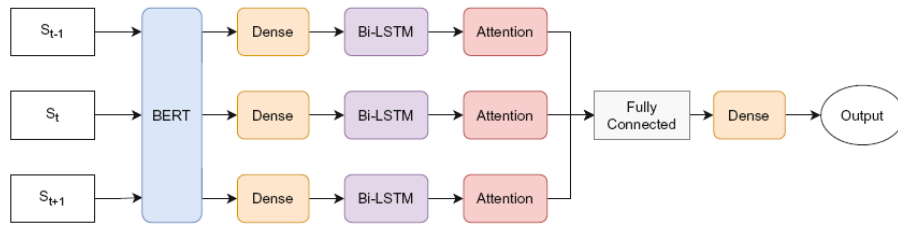


Fig. 1. Model architecture. One BERT structure is used to encode each sentence.

1. BERT is used as initial sentence encoder.
2. First hidden dense layer acts as a feature extraction layer, mapping the highly dimensional sentence embedding to a vector with a smaller dimension.
3. Bi-directional LSTM layer captures longer relational structure within the sentence, only takes 3 sentences (previous, current, next) into account.
4. Attention layers capture the important pieces of the sentences.
5. Final dense layer takes the 3 independent streams/sentence encodings and produces an output.

Dataset consists of book chapters and annotated Wikipedia articles. Data for training was augmented by introducing a max segment length of 5 (or less), reducing the imbalance between the positive and negative class. This did come with problems, the significantly reduced dataset size made overfitting easier.

Bibliography

- [1] S. Arnold, R. Schneider, P. Cudré-Mauroux, F. A. Gers and A. Löser, "SECTOR: A Neural Model for Coherent Topic Segmentation and Classification," *Transactions of the Association for Computational Linguistics*, pp. 169-184. doi: https://doi.org/10.1162/tacl_a_00261, 2019.
- [2] O. Koshorek, A. Cohen, N. Mor, M. Rotman and J. Berant, "Text Segmentation as a Supervised Learning Task," *Proceedings of NAACL-HLT*, pp. 469-473. doi: <https://doi.org/10.48550/arXiv.1803.09337>, 2018.
- [3] I. Sheikh, D. Fohr and I. Illina, "Topic segmentation in ASR transcripts using bidirectional rnns for change detection," *2017 IEEE automatic speech recognition and understanding workshop (ASRU)*, pp. 512-518, doi: <https://doi.org/10.1109/ASRU.2017.8268979>, 2017.
- [4] A. Maraj, M. V. Martin and M. Makrehchi, "A More Effective Sentence-Wise Text Segmentation Approach Using BERT," *International Conference on Document Analysis and Recognition*, pp. 236-250, doi: https://doi.org/10.1007/978-3-030-86337-1_16, 2021.

