

Content Selection in Deep Learning Models of Summarization

Anonymous EMNLP submission

Abstract

1 Introduction

While there has been a recent flurry of work on abstractive summarization (????), these papers treat this problem as a pure sequence to sequence transduction task. Admittedly, this view allows us to apply very powerful, general-purpose deep learning architectures to generate summaries. At the same time, it obscures a principal subtask in summarization, the process of selecting the most salient units of meaning in the source material, i.e. the key ingredients in the final summary, a process which we broadly refer to as content selection (?).

As is also the case in other NLP tasks, it is not immediately obvious how a deep learning model is making its predictions, or what correlations are being exploited. There is a concerning and growing list of papers that find models functioning as mere nearest neighbors search (??), exploiting annotator artifacts (?), or open to adversarial exploitation (?). These lines of research are critical for finding model shortcomings, and over time, guiding improvements in technique. Unfortunately, to the best of our knowledge, there has been no such undertaking for the summarization task.

In this paper, we seek to better understand how deep learning models of summarization are performing content selection. We perform an analysis of several recent sentence extractive neural network architectures, looking particularly at the impact of sentence position bias, the necessity of learning embeddings, the unreasonable effectiveness of averaging for sentence embedding, and the cross domain generalizability of such models. Additionally, we propose two simpler models that are on average statistically indistinguishable from their more complex counterparts.

While we are explicitly studying extractive summarization algorithms here, we think the findings will be relevant to the abstractive summarization community as well. The encoder side architectures are quite similar to typical abstractive models, and fundamentally the model objectives are the same, producing output text with high word overlap to a reference human abstract.

The contributions of this paper are the following:

1. We perform an empirical study of extractive content selection in deep learning algorithms for text summarization across news, blogs, meetings, and journal articles domains, and small/medium/large datasets.
2. Propose two simple deep learning models whose performance is on par with more complex deep learning models.

In the following sections we discuss (Sec. ?) related work, (Sec. ?) define the problem of extractive summarization, (sec. ?) formulate our proposed ad baseline models, (Sec. ?) describe the datasets used for experiments, (Sec. ?) describe the experiments themselves, and conclude with results and analysis (Sec. ?).

2 Related Work

Extractive Deep Learning Based Summarization

Nallapati et. al
Cheng and Lapata

Abstractive Deep Learning Based Summarization

Rush
Chopra
See et al.

Socher

Extractive Single Doc Summarization Durrett
et. al

Non Newswire Summarization meeting
summarization
reddit stories
journal articles/pubmed

3 Problem Definition

The goal of extractive text summarization is to select a subset of a document’s text to use as a summary, i.e. a short gist or excerpt of the central content. Typically, we impose a budget on the length of the summary in either words or bytes.

In this paper, we model this task as a sequence tagging problem, i.e. given a document containing d sentences s_1, \dots, s_d we want to predict a corresponding label sequence $y_1, \dots, y_d \in \{0, 1\}^d$ where $y_i = 1$ indicates the i -th sentence is to be included in the summary.

Unlike sequence tagging, however, we do not evaluate model performance by label accuracy or F-measure but ROUGE (?) which measures the ngram overlap of our predicted extract summary with one or more human abstracts. While this metric has many shortcomings, ROUGE-2 recall (i.e. bigram recall) has been empirically demonstrated to have high correlation with human content selection decisions for summarization (?).

Since we do not typically have ground truth extract summaries from which to create the labels y_i , we construct gold label sequences for training by greedily optimizing ROUGE-1. Starting with an empty summary $S = \emptyset$, we add the sentence $\hat{s} = \arg \max_{s \in \{s_1, \dots, s_d\}, s \notin S} \text{ROUGE-1}(S \cup s)$ to S stopping when the ROUGE-1 score no longer increases or the length budget is reached. We choose to optimize for ROUGE-1 rather than ROUGE-2 similarly to other optimization based approaches to summarization (????) which found this be the easier optimization target.

4 Models and Methods

At a high level, all of our models share the same two part structure: i) a *sentence encoder* which maps an arbitrary sequence of tokens to an embedding $h \in \mathcal{R}^d$, and ii) a *sentence extractor* which takes as input all of a document’s sentence em-

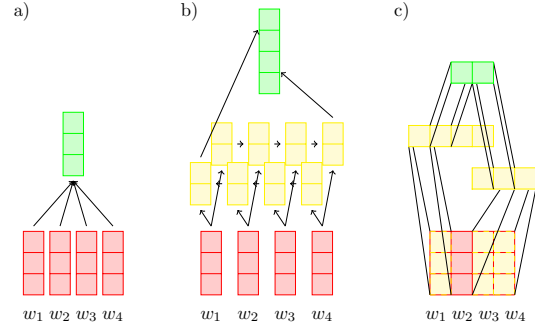


Figure 1: Sentence encoder architectures: a) averaging encoder, b) RNN encoder c) CNN encoder. Red indicates word embeddings, yellow indicates RNN hidden states or convolutional activations, and green indicates the sentence embedding that is passed to the extractor module.

beddings and predicts which sentences to extract to produce the extract summary. The sentence extractor is essentially a discriminative classifier $p(y_1, \dots, y_d | h_1, \dots, h_d)$.

Depending on the architectural choices of each component we propose we can recover the specific implementations of (?) and (?), which we outline below.

4.1 Sentence Encoders

We treat each sentence $s = \{w_1, \dots, w_{|s|}\}$ as a sequence of word embeddings, where $|s|$ is the total number of words in the sentence. We experiment with three architectures for mapping sequences of word embeddings to a fixed length vector: average pooling, RNNs, and CNNs.

Average Pooling Average pooling is the simplest method and parameter free. A sentence encoding $\text{enc}(s) = \frac{1}{|s|} \sum_{w \in s} w$.

RNN Encoder Our second sentence encoder uses the concatenation of the final output states of a forward and backward RNN over the sentence’s word embeddings. We use a Gated Recurrent Unit (GRU) (?) for the RNN cell, since it has fewer parameters than the equivalent LSTM but with similar performance.

$$\vec{h}_i = \overrightarrow{\text{GRU}}(w_i, h_{i-1}) \quad \forall i \in 1, \dots, |s| \quad (1)$$

$$\overleftarrow{h}_i = \overleftarrow{\text{GRU}}(w_i, h_{i+1}) \quad \forall i \in 1, \dots, |s| \quad (2)$$

$$\text{enc}(s) = [\vec{h}_{|s|}; \overleftarrow{h}_1] \quad (3)$$

CNN Encoder Our final sentence encoder uses a series of convolutional feature maps to encode

each sentence. This encoder is similar to the convolutional architecture of (?) used for text classification tasks and performs a series of “one-dimensional” convolutions over a sentence’s associated word embeddings. In addition to its learned parameters, the CNN encoder has hyperparameters associated to the window size of the convolutional filter (i.e. the number of words associated with each convolution) and the number of feature maps associated with each filter (i.e. the output dimension of each convolution). The CNN sentence encoding is computed as follows:

$$a_i^{(f,m)} = b^{(f,m)} + \sum_{j=1}^m W_j^{(f,m)} \cdot w_{i+j-1} \quad (4)$$

$$h_i^{(f,m)} = \max_{i \in 1, \dots, |s|-m+1} \text{ReLU} \left(a_i^{(f,m)} \right) \quad (5)$$

$$\text{enc}(s) = \left[h^{(f,m)} \mid \forall f \in F, m \in M \right] \quad (6)$$

4.2 Sentence Extractors

Given a sequence of sentence embeddings $h_i = \text{enc}(s_i)$, a sentence extractor produces a conditional distribution over the corresponding sentence extraction variables $p(y_1, \dots, y_{|s|} \mid h_1, \dots, h_{|s|})$. We propose two simple recurrent neural network based sentence extractors that make a strong conditional independence assumption over the labels y_i , chiefly $p(y_1, \dots, y_{|s|} \mid h_1, \dots, h_{|s|}) = \prod_{i=1}^{|s|} p(y_i \mid h_1, \dots, h_{|s|})$. This stands in contrast to our baseline models which make a weaker assumption, $p(y \mid h) = \prod_{i=1}^{|s|} p(y_i \mid y_{<i}, h_1, \dots, h_{|s|})$, at the expense of greater computational complexity.

RNN Extractor Our first proposed model is a very simple bidirectional RNN based tagging model. As in the RNN sentence encoder we use a GRU cell. The the forward and backward outputs of each sentence are passed through a single layer perceptron with a sigmoid output to predict the probability of extracting each sentence:

$$\vec{z}_i = \overrightarrow{\text{GRU}}(h_i, \vec{z}_{i-1}) \quad (7)$$

$$\overleftarrow{z}_i = \overleftarrow{\text{GRU}}(h_i, \overleftarrow{z}_{i+1}) \quad (8)$$

$$a_i = \text{ReLU} (U \cdot [\vec{z}_i; \overleftarrow{z}_i] + u) \quad (9)$$

$$p(y_i = 1 \mid h) = \sigma (V \cdot a_i + v). \quad (10)$$

Seq2Seq Extractor One shortcoming of the RNN extractor is that long range information from one end of the document may not easily be able

to effect extraction probabilities of sentences at the other end. To mitigate this effect we introduce a Seq2Seq extractor based on the attentional seq2seq models commonly used for neural machine translation (?) and abstractive summarization (?). The sentence embeddings are first encoded by a bidirectional GRU. A separate decoder GRU transforms each sentence into a query vector which attends to the encoder output. The attention weighted encoder output and the decoder GRU output are concatenated and fed into a multi-layer perceptron to compute the extraction probability. Formally we have:

$$\vec{z}_i = \overrightarrow{\text{GRU}}_{\text{enc}}(h_i, \vec{z}_{i-1}) \quad (11)$$

$$\overleftarrow{z}_i = \overleftarrow{\text{GRU}}_{\text{enc}}(h_i, \overleftarrow{z}_{i+1}) \quad (12)$$

$$\vec{q}_i = \overrightarrow{\text{GRU}}_{\text{dec}}(h_i, \vec{q}_{i-1}) \quad (13)$$

$$\overleftarrow{q}_i = \overleftarrow{\text{GRU}}_{\text{dec}}(h_i, \overleftarrow{q}_{i+1}) \quad (14)$$

$$\alpha_{i,j} = \frac{\exp(\vec{q}_j \cdot \vec{z}_i + \overleftarrow{q}_j \cdot \overleftarrow{z}_i)}{\sum_{j=1}^d \exp(\vec{q}_j \cdot \vec{z}_i + \overleftarrow{q}_j \cdot \overleftarrow{z}_i)} \quad (15)$$

$$\bar{z}_i = \sum_{j=1}^d \alpha_{i,j} [\vec{z}_j; \overleftarrow{z}_j] \quad (16)$$

$$a_i = \text{ReLU} (U \cdot [\bar{z}_i; \vec{q}_i; \overleftarrow{q}_i] + u) \quad (17)$$

$$p(y_i = 1 \mid h) = \sigma (V \cdot a_i + v). \quad (18)$$

Cheng & Lapata Extractor

$$z_i = \text{GRU}_{\text{enc}}(h_i, z_{i-1}) \quad (19)$$

$$q_i = \text{GRU}_{\text{dec}}(p_{i-1} \cdot h_{i-1}, q_{i-1}) \quad (20)$$

$$a_i = \text{ReLU} (U \cdot [z_i; q_i] + u) \quad (21)$$

$$p_i = p(y_i = 1 \mid y_{<i}, h) = \sigma (V \cdot a_i + v) \quad (22)$$

SummaRunner Extractor

$$\vec{z}_i = \overrightarrow{\text{GRU}}(h_i, \vec{z}_{i-1}) \quad (23)$$

$$\overleftarrow{z}_i = \overleftarrow{\text{GRU}}(h_i, \overleftarrow{z}_{i+1}) \quad (24)$$

$$q = \tanh \left(b_q + W_q \frac{1}{d} \sum_{i=1}^d [\vec{z}_i; \overleftarrow{z}_i] \right) \quad (25)$$

$$z_i = \text{ReLU} (b_z + W_z [\vec{z}_i; \overleftarrow{z}_i]) \quad (26)$$

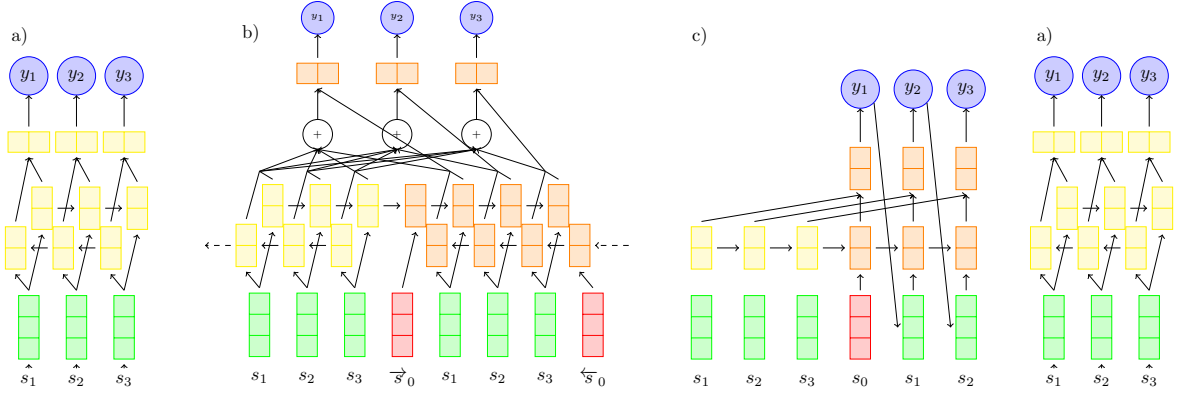


Figure 2: Sentence extractor architectures: a) RNN, b) Seq2Seq, c) C&L, and d) SR.

Figure 3: fig:extractor

$$\begin{aligned}
 p(y_i = 1 | y_{<i}, h) = & \sigma(W_{con} \cdot z_i \\
 & + z_i^T W_{sal} \cdot q \\
 & - z_i^T W_{nov} \cdot \tanh(g_i) \\
 & + b_{rp_i} \\
 & + b_{ap_i} \\
 & + b) \\
 g_j = & \sum_{i=1}^{j-1} p(y_j = 1 | y_{<j}, h) \cdot z_j
 \end{aligned}$$

5 Datasets

cnn-dailymail
 nyt
 duc
 ami
 reddit
 pubmed

6 Experiments

Extractor/Encoder In our main experiment we compare our proposed sentence extractors **RNN** and **Seq2Seq** against the **C&L** and **SR** extractors. We test all possible sentence extractor/encoder pairs across the CNN-DailyMail, New York Times, DUC 2002, Reddit, AMI, and PubMed domains. We choose ROUGE-2 recall as our main evaluation metric since it has the strongest correlation to human content selection decisions. In this experiment we initialize the word embeddings using pretrained GloVe embeddings (?) and do not update them during training.

In most cases, the averaging encoder performance was as good or better than the RNN and CNN encoders, we use only the averaging encoder for the remainder of the experiments.

Word Embedding Learning To further understand how word embeddings can effect model performance we also compared extractors when embeddings are updated during training. Both fixed and learned embedding variants are initialized with GloVe embeddings. When learning embeddings, words occurring three or fewer times in the training data are mapped to an *unknown* token.

POS Tag Ablation Additionally, we ran ablation experiments using part-of-speech (POS) tags. We experimented with selectively removing nouns, verbs, adjectives/adverbs, numerical expressions, and miscellaneous tags (anything that was not in the previously mentioned groups) from each sentence. The embeddings of removed words were replaced with a zero vector, preserving the order and position of the non-ablated words in the sentence. All datasets were automatically tagged using the SpaCy POS tagger (?).

Document Shuffling In examining the outputs of the models, we found most of the selected sentences in the news domain came from the lead paragraph of the document. This is despite the fact that there is a long tail of sentence extractions from later in the document in the ground truth extract summaries. Because this lead bias is so strong, it is questionable whether the models are learning to identify important content or just find the start of the document. We perform a series of sentence order experiments where each document's sentences

are randomly shuffled during training. We then evaluate each model performance on the unshuffled test data, comparing to the original unshuffled models.

Cross Domain Experiments

TODO

Try shuffled and no shuffled models trained on one domain, eval the remaining.

6.1 Model Training Details

TODO: Clean up, expand, make naming consistent!

We train the average pooling, biRNN, and CNN encoders with the biRNN, seq2seq, SummaRunner, and C&L decoders. We repeat experiments across the CNN-DailyMail, New York Times, DUC, Reddit, and AMI corpus. We use the Adam optimizer for all models with a learning rate of .0001, gradient clipping, and dropout rate of .25. For the C&L model, we train for half of the maximum epochs with teacher forcing, i.e. we use the gold extractive labels for each when taking the sum of previous states ($p(y_{t-1}|x) = 1$ if $y_t = 1$) during the first half of training and the model value during the second. We use early stopping on the validation set with ROUGE2 as our evaluation criteria. All experiments are run with 5 different random initialization seeds and results are averaged.

7 Results

Choice of Extractor Our main comparison of extractors/encoders are shown in Table 1. Overall we find that the Seq2Seq extractor achieves the best ROUGE scores on three out of four domains (STILL RUNNING ON AMI AND PUBMED). However, most differences are not significant. (Need to discuss stat sig and how to show it). On the larger CNN-DailyMail dataset, especially, differences are quite small across all extractor/encoder pairs. The C&L extractor achieves the best performance on the DUC 2002 dataset. It is disappointing that the C&L and SR based models do not gain any apparent advantage in conditioning on previous sentence selection decisions; this result suggests the need to improve the representation of the summary as it is being constructed iteratively.

Choice of Encoder We also find there to be no major advantage between the different sentence encoders. In most cases, there is no statistical significance between the averaging encoder and ei-

ther the RNN or CNN encoders.

Learning Word Embeddings Table 2 shows ROUGE recall when using fixed or updated word embeddings. In almost all cases, fixed embeddings are as good or better than the learned embeddings.

POS Ablation Table 3 shows the results of the POS tag ablation experiments. The newswire domain does not appear to be sensitive to these ablations; this suggests that the models are still able to identify the lead section of the document with the remaining word classes (Verify this with histogram analysis). The Reddit domain, which is not lead biased, is significantly effected. Notably, removing adjectives and adverbs results in a 1.8 point drop in ROUGE-2 recall.

Sentence Shuffling We find a similar result on the sentence order shuffling experiments. Table 4 shows the results. The newswire domain suffer a significant drop in performance when the document order is shuffled. By comparison, there is no significant difference between the shuffled and in-order models on the Reddit domain.

| Extractor | Encoder | cnn-dailymail | | nyt | | duc-sds | | reddit | |
|-----------|---------|---------------|-------------|------|-------------|-------------|-------------|-------------|-------------|
| | | R-1 | R-2 | R-1 | R-2 | R-1 | R-2 | R-1 | R-2 |
| RNN | avg | 55.3 | 25.4 | 51.4 | 34.7 | 44.1 | 22.6 | 45.2 | 11.4 |
| | rnn | 55.3 | 25.4 | 51.7 | 35.0 | 44.0 | 22.6 | 44.4 | 11.4 |
| | cnn | 54.7 | 25.1 | 50.3 | 33.7 | 44.3 | 22.7 | 47.5 | 12.7 |
| Seq2Seq | avg | 55.6 | 25.6 | 52.5 | 35.7 | 44.4 | 22.8 | 49.1 | 13.6 |
| | rnn | 55.2 | 25.3 | 52.5 | 35.9 | 43.8 | 22.5 | 45.4 | 12.1 |
| | cnn | 54.8 | 25.1 | 51.7 | 35.1 | 44.0 | 22.7 | 46.8 | 13.1 |
| C&L | avg | 55.1 | 25.3 | 52.3 | 35.6 | 44.8 | 23.1 | 48.3 | 13.6 |
| | rnn | 54.8 | 25.0 | 52.5 | 35.8 | 44.5 | 23.0 | 46.3 | 12.6 |
| | cnn | 55.0 | 25.1 | 51.5 | 35.0 | 44.5 | 23.0 | 46.9 | 13.4 |
| SR | avg | 55.3 | 25.4 | 52.1 | 35.4 | 44.0 | 22.3 | 48.8 | 13.4 |
| | rnn | 55.0 | 25.2 | 52.1 | 35.5 | 43.6 | 22.1 | 46.5 | 12.6 |
| | cnn | 54.6 | 25.0 | 51.0 | 34.4 | 43.6 | 22.2 | 46.6 | 12.3 |

Table 1: Rouge Recall 1 and 2 results across all sentence encoder/extractor pairs. All results are averaged over five random initializations. Best result per metric/dataset are bolded.

| system | embeddings | cnn-dailymail | | nyt | | duc-sds | | reddit | |
|-------------|------------|---------------|------|------|------|---------|------|--------|------|
| | | R1 | R2 | R1 | R2 | R1 | R2 | R1 | R2 |
| RNN | fixed | 55.3 | 25.4 | 51.4 | 34.7 | 44.1 | 22.6 | 45.2 | 11.4 |
| | learned | 55.1 | 25.2 | 51.1 | 34.3 | 44.1 | 22.6 | 45.3 | 11.3 |
| Seq2Seq | fixed | 55.6 | 25.6 | 52.5 | 35.7 | 44.4 | 22.8 | 49.1 | 13.6 |
| | learned | 55.2 | 25.3 | 52.4 | 35.7 | 44.5 | 22.9 | 49.4 | 13.8 |
| C&L | fixed | 55.1 | 25.3 | 52.3 | 35.6 | 44.8 | 23.1 | 48.3 | 13.6 |
| | learned | 54.8 | 25.0 | 52.1 | 35.4 | 44.6 | 23.0 | 48.6 | 13.5 |
| SummaRunner | fixed | 55.3 | 25.4 | 52.1 | 35.4 | 44.0 | 22.3 | 48.8 | 13.4 |
| | learned | 55.0 | 25.1 | 52.0 | 35.2 | 43.8 | 22.1 | 47.8 | 12.6 |

Table 2: ROUGE 1 and 2 recall results across different sentence extractors when using learned or pretrained embeddings. In both cases embeddings are initialized with pretrained GloVe embeddings. All results are averaged from five random initializations. All extractors use the averaging sentence encoder.

| Ablation | cnn-dailymail | | nyt | | duc-sds | | reddit | |
|----------|---------------|------|------|------|---------|------|--------|------|
| | R1 | R2 | R1 | R2 | R1 | R2 | R1 | R2 |
| – | 55.3 | 25.4 | 51.4 | 34.7 | 44.1 | 22.6 | 45.2 | 11.4 |
| -nouns | 55.2 | 25.3 | 50.8 | 34.3 | 43.8 | 22.3 | 43.0 | 10.3 |
| -verbs | 55.1 | 25.3 | 51.0 | 34.4 | 43.8 | 22.4 | 44.3 | 10.7 |
| -adjv | 55.2 | 25.3 | 51.0 | 34.4 | 44.0 | 22.5 | 42.5 | 9.6 |
| -misc | 55.1 | 25.2 | 51.2 | 34.5 | 44.5 | 22.9 | 43.1 | 10.3 |
| -num | 55.3 | 25.4 | 51.3 | 34.6 | 44.1 | 22.6 | 45.2 | 11.0 |

Table 3: ROUGE recall after removing different word classes. Ablations are performed using the averaging sentence encoder and **RNN** extractor. Table shows average results of five random initializations.

| Extractor | sentence order | nyt | | duc-sds | | reddit | |
|-----------|----------------|------|------|---------|------|--------|------|
| | | R1 | R2 | R1 | R2 | R1 | R2 |
| RNN | in-order | 51.4 | 34.7 | 44.1 | 22.6 | 45.2 | 11.4 |
| | shuffled | 41.9 | 25.0 | 39.7 | 18.2 | 45.1 | 11.9 |
| Seq2Seq | in-order | 52.5 | 35.7 | 44.4 | 22.8 | 49.1 | 13.6 |
| | shuffled | 42.6 | 25.6 | 42.9 | 21.2 | 48.7 | 13.6 |

Table 4: ROUGE 1 and 2 recall using models trained on in-order and shuffled documents. All extractors use the averaging sentence encoder. Table shows average results of five random initializations.