

Deep Q-Learning for Event Summarization

Francisco Javier Arceo

fja2114@columbia.edu

Chris Kedzie

kedzie@cs.columbia.edu

November 12, 2016

Abstract

We present a streaming extraction policy for multi-document summarization learned using a deep reinforcement learning recurrent neural network architecture. The rewards of our network are evaluated through manually reviewed summaries, which allow us to specify an end-to-end framework to optimize our extraction policy. By using raw text from our articles and mapping words into a higher embedding dimension through our recurrent neural network, we are able to learn a more robust representation than traditional n-gram models. We benchmark our model against random actions, a bag-of-words model, and a bag-of-bigrams model.¹

1 Introduction

Recent research in text retrieval has focused on extractive algorithms to identify important sentences from a large set of documents (e.g., [?], [?], [?], and [?]). The approach by [?] has shown that it is possible to select relevant sentences from a massive number of documents on the web to create summaries with meaningful content by adapting classifiers to maximize search policies. These systems operate in a streaming fashion and are capable of evaluating each sentence within each article to decide whether or not to include or ignore the sentence.

Unfortunately, many of these systems have still been shown to fall short of algorithms that employ simple heuristics [?], which may be due to inadequate capturing of the rich structure and often idiosyncratic information by traditional n-gram language models.

¹ <https://github.com/franciscojavierarceo/DQN-Event-Summarization>

In recent years, recurrent neural networks have been used to map textual language to a higher order dimension to represent more powerful features for a variety of different language modeling tasks [?]. These vector representations, also known as embeddings, have proven incredibly powerful on a variety of different natural language processing tasks (e.g., [?], [?], [?]). For query based retrievals, [?] has shown success in using these embeddings for *ad hoc* information retrieval.

In this paper we show that a deep recurrent neural network with a long short term memory (LSTM) [?] is able to successfully learn an extractive summary policy by encoding the action, state, and reward into a Q-Learning model, similar to that of [?]. We show the performance of our model on a variety of different metrics and benchmark them against random actions, a bag-of-words multilayered perceptron, and bag-of-bigrams multilayered perceptron.

2 Extractive Streaming Summarization

In the extractive streaming summarization task, we are given as input a query (i.e., a short text description of a topic or an event), a document stream, and a time ordered set of sentences relevant to the query. Starting with an initially empty summary, an extractive, streaming summarization algorithm is intended to examine each sentence in order and, when new and important (relative to the query) information is identified, add that sentence to the summary.

Implicit to this problem is the notion of system time – the summarization algorithm can only examine sentences that occur in the stream before the current system time. Advancing the system time gives the algorithm access to more sentences, although in practice the stream is sufficiently large enough that choices have to be made about how much history can be kept in memory. For many domains, e.g., crisis informatics, it is preferable for a summarization algorithm to identify important information as early as possible, and so the objective function should penalize a large disparity between the time a piece of information is first available to the algorithm and the system time at which that information is actually added to the summary.

Previous work in this area has either incremented the system time in fixed increments (e.g., an hour) [?, ?] or operated in a fully online setting [?, ?]. In both cases explicitly modeling the current state of the summary, the stream, and their relationship with the query is quite complicated and exhibits non-linear dynamics that are difficult to characterize in traditional feature based models.

Additionally, the structured nature of the sentence selection task (sentence selection is highly dependent on the current summary state) suggests that imitation or

reinforcement learning are necessary to obtain parity between training and testing feature distributions.

Early successes of deep reinforcement learning have used convolutional neural networks for video games (e.g., [?] and [?]) but recent work has leveraged recurrent neural networks for both video-based and text-based gaming environments (e.g., [?] and [?]).

This leads us to explore Deep Q-Networks (DQN) for three reasons: (1) both the representation and interaction between the stream, summary, and query can be learned; (2) the embeddings can learn a more robust semantic representations than classic n-gram models; and (3) by randomly exploring the state space, the ϵ -greedy strategy used in DQN learns a policy that yields more consistency between train and test distributions.

3 Background

In the subsequent sections, we outline the architecture of our DQN and define the parameterization of our Q-function for the multi-document summarization task. Our architecture consists of three components: the (1) input, (2) action, and (3) value representations. By specifying these components into our neural network, we are able to develop an end-to-end extraction policy that fully propagates information both forward and backward.

3.1 States

In our architecture a state $s(x_t, \tilde{y}_t, d)$ is a function of the stream X_d , $\forall d \in D$, the state of the current summary \tilde{y}_t at system time t , and the query d . For brevity we will use s_t where the dependence on x_t, \tilde{y}_t , and d is assumed. s_t itself is three recurrent neural networks, one for encoding the summary, the stream, and the query.

3.2 Actions

The set of possible actions a_t at each time step is $\mathcal{A} = \{select, skip\}$ where *select* corresponds to adding the current sentence x_t to the summary and incrementing the current system time, or *skip* where only t is incremented without changing the current summary. Thus, our current predicted summary at time t is defined as

$$\tilde{y}_{t+1} = \begin{cases} \tilde{y}_t \cup \{x_t\}, & \text{if } a_t = select \\ \tilde{y}_t, & \text{if } a_t = skip. \end{cases} \quad (1)$$

Since \tilde{y}_t grows linearly with the number of selected sentences, we mitigate potential capacity constraints by using a queue to hold the last K elements of the current predicted summary.

3.3 Reward

The reward for a given action is measured by the change in ROUGE-N F1 score of the predicted summary \tilde{y}_t measured against a gold standard summary Y . When only one gold summary reference is used, ROUGE-N Recall is calculated as

$$\text{ROUGE-NR}(\tilde{y}, Y) = \frac{\sum_{g \in \text{ngrams}(Y, N)} \min(\text{count}(g, \tilde{y}), \text{count}(g, Y))}{\sum_{g \in \text{ngrams}(Y, N)} \text{count}(g, Y)} \quad (2)$$

where $\text{ngrams}(Y, N)$ returns the set of ngrams of order N in the summary Y and $\text{count}(g, Y)$ is the count of occurrences of ngram g in Y .

Similarly, ROUGE-N Precision is calculated as

$$\text{ROUGE-NP}(\tilde{Y}, Y) = \frac{\sum_{g \in \text{ngrams}(Y, N)} \min(\text{count}(g, \tilde{Y}), \text{count}(g, Y))}{\sum_{g \in \text{ngrams}(\tilde{Y}, N)} \text{count}(g, \tilde{Y})} \quad (3)$$

and the F_1 is simply the harmonic mean of the two:

$$\text{ROUGE-NF1}(\tilde{Y}, Y) = \frac{2 \times \text{ROUGE-NP}(\tilde{Y}, Y) \times \text{ROUGE-NR}(\tilde{Y}, Y)}{\text{ROUGE-NP}(\tilde{Y}, Y) + \text{ROUGE-NR}(\tilde{Y}, Y)} \quad (4)$$

The reward r at time t is

$$r_t = \text{ROUGE-NF1}(\tilde{y}_t, Y) - \text{ROUGE-NF1}(\tilde{y}_{t-1}, Y). \quad (5)$$

3.4 Deep Q-Learning

We define an architecture similar to that of [?] and map our three inputs (query, sentence, and current predicted summary) into LSTM embeddings according to **Figure 1**. By formulating our extraction task as a Markov Decision Process (MDP) we can express the state-action tuples as transitions states.

Our extraction policy then takes as input an action a_t in state s_t at time t and returns the expected reward. We initialize our actions and Q-function at random and by taking the optimal action given our expected reward, we are able to iteratively update our Q-function using the Bellman equation [?] satisfying

$$\hat{Q}_{t+1}(s, a) = \mathbb{E}[r + \gamma \max_{a'} \hat{Q}_t(s', a') | s, a] \quad (6)$$

where $\gamma \in \mathbb{R}$ is the discount rate applied to future rewards and the expectation is taken over transitions in state s and action a . By iterating over the Q-function, it converges asymptotically to the optimal policy, i.e., $\hat{Q}_t \rightarrow Q^*$ as $t \rightarrow \infty$, but in finite iterations we can use a neural network as our function approximator, where our weights, θ , can be optimized using the following loss function

$$\mathcal{L}_i(\theta_i) = \mathbb{E}_{\hat{s}, \hat{a}} [(y_i - Q(\hat{s}, \hat{a}; \theta_i))^2] \quad (7)$$

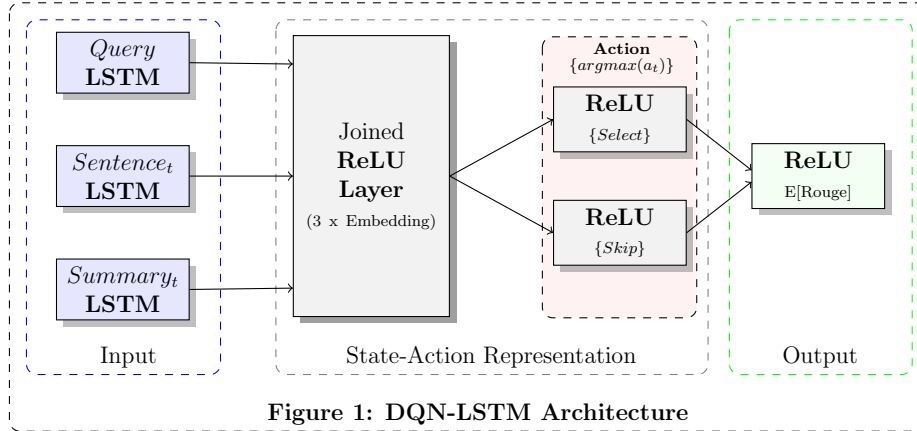
where

$$y_i = \mathbb{E}_{\hat{s}, \hat{a}} [r_i + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s', a']. \quad (8)$$

Differentiation of the loss function yields the update equation,

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{\hat{s}, \hat{a}} [2(y_i - Q(\hat{s}, \hat{a}; \theta_i)) \nabla_{\theta_i} Q(\hat{s}, \hat{a}; \theta_i)] \quad (9)$$

that learns the parameters using stochastic gradient descent with RMSprop [?] on sampled batches of stored transitions.



4 Learning an Extraction Policy

4.1 Settings

Since the state of our current summary at time t depends on the state at time $t - 1$, we have to execute our Q-function in the feedforward stage using a single batch to

properly update the predicted summary with optimal expected action. During the backpropagation phase we sample random mini-batches from the stored transitions. We set our embedding dimension to 100, the future discount rate γ to 0.50, and the maximum summary length to 300 tokens.

The weights of the Q-function and the action sequence are initialized at random and we randomly explore the state-space for 200 of the 1000 epochs by employing an ϵ -greedy search strategy during our action selection that linearly decays to a base exploration rate of 0.01. We describe in full detail the training of our DQN-LSTM in **Algorithm 1**.

Algorithm 1 DQN-LSTM for Event Summarization Training Procedure

Input: $\{\mathcal{D}$: Event queries, X_d : Input sentences, N : Number of epochs $\}$

Output : $\{\hat{Q}$: extraction policy, \tilde{Y}_d : event summary for query $d\}$

- 1: Initialize action-value function \hat{Q} with random weights
 - 2: Initialize memory and summary: $\Gamma, \tilde{Y} = \{\emptyset\}_{d=1}^{|\mathcal{D}|}, \{\emptyset\}_{d=1}^{|\mathcal{D}|}$
 - 3: **for** $epoch = 1, \dots, N$ **do**
 - 4: **for** query $d \in \mathcal{D}$ **do**
 - 5: $X_d, \tilde{Y}_d = \{\text{Extract } t = 1, \dots, T_d \text{ (sentences}_d, \text{summary}_d)\}$
 - 6: **for** $x_t, \tilde{y}_t \in X_d, \tilde{Y}_d$ **do**
 - 7: Set $s_t = s(x_t, \tilde{y}_t, d)$
 - 8: $\forall a_t \in \mathcal{A}(s_t)$ compute $\hat{Q}(s_t, a_t)$ and select $a_t^* = \text{argmax}_{a_t} \hat{Q}(s_t, a_t)$
 - 9: **if** $\text{random}() < \epsilon$ **then** select a_t^* at random with $\text{Pr}(a_t) = \frac{1}{|\mathcal{A}|}$
 - 10: Update \tilde{y}_{t+1} according to equation (1)
 - 11: Execute action a_t^* and observe reward r_t and new state s_{t+1}
 - 12: Update $\Gamma_d = \Gamma_d \cup \{[s_t, a_t^*, r_t, s_{t+1}]\}$
 - 13: **for** $j = 1, \dots, J$ transitions sampled from Γ **do**
 - 14: Set $y_j = \begin{cases} r_j & \text{if } s_{j+1} \text{ is terminal} \\ r_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a'; \theta) & \text{if } s_{j+1} \text{ is non-terminal} \end{cases}$
 - 15: Perform gradient step on $\mathcal{L}(\theta) = (y_j - \hat{Q}(s_j, a_j; \theta))^2$
-

4.2 Simulations

To understand the behavior of our architecture, we produce a number of experiments on the first sentence of the first 100 articles for a single query. By restricting our

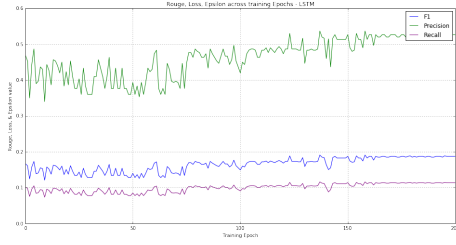


Figure 1: Image 1

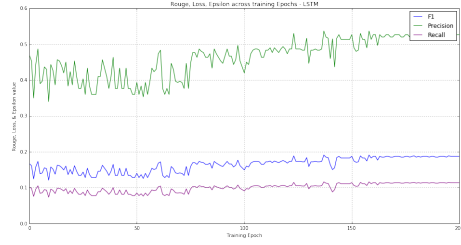


Figure 2: Image 2

Figure 2a: DQN-LSTM Hyperparameter Experiments

study to such small data we are able to better monitor the system and influence of hyperparameters. We study the performance of hyperparameters on a small set of test data in the next section.

In **Figure 2a** we see the convergence of our model when trained on F1, which suggests that our Q-function learns a policy that tries to vary the Precision while maintaining a fixed Recall. When we specify the model to learn Recall or Precision obvious pathological policies are learned, i.e., our model decides to select all sentences or select only the fewest that maximize the intersection with the n-grams in our nuggets. In **Figure 2b** we present the model without a base exploration rate, which suggests that adding a small amount of randomness helps the model from converging to poor local optimums.