# Recurrent Neural Network for Text Classification
# with Multi-Task Learning

**Pengfei Liu    Xipeng Qiu*    Xuanjing Huang**

Shanghai Key Laboratory of Intelligent Information Processing, Fudan University

School of Computer Science, Fudan University

825 Zhangheng Road, Shanghai, China

{pfliu14,xpqiu,xjhuang}@fudan.edu.cn

## Abstract

Neural network based methods have obtained great progress on a variety of natural language processing tasks. However, in most previous works, the models are learned based on single-task supervised objectives, which often suffer from insufficient training data. In this paper, we use the multi-task learning framework to jointly learn across multiple related tasks. Based on recurrent neural network, we propose three different mechanisms of sharing information to model text with task-specific and shared layers. The entire network is trained jointly on all these tasks. Experiments on four benchmark text classification tasks show that our proposed models can improve the performance of a task with the help of other related tasks.

## 1 Introduction

Distributed representations of words have been widely used in many natural language processing (NLP) tasks. Following this success, it is rising a substantial interest to learn the distributed representations of the continuous words, such as phrases, sentences, paragraphs and documents [Socher *et al.*, 2013; Le and Mikolov, 2014; Kalchbrenner *et al.*, 2014; Liu *et al.*, 2015a]. The primary role of these models is to represent the variable-length sentence or document as a fixed-length vector. A good representation of the variable-length text should fully capture the semantics of natural language.

The deep neural networks (DNN) based methods usually need a large-scale corpus due to the large number of parameters, it is hard to train a network that generalizes well with limited data. However, the costs are extremely expensive to build the large scale resources for some NLP tasks. To deal with this problem, these models often involve an unsupervised pre-training phase. The final model is fine-tuned with respect to a supervised training criterion with a gradient based optimization. Recent studies have demonstrated significant accuracy gains in several NLP tasks [Collobert *et al.*, 2011] with the help of the word representations learned from the large unannotated corpora. Most pre-training methods are based on unsupervised objectives such as word prediction for training [Collobert *et al.*, 2011; Turian *et al.*, 2010; Mikolov *et al.*, 2013]. This unsupervised pre-training is effective to improve the final performance, but it does not directly optimize the desired task.

Multi-task learning utilizes the correlation between related tasks to improve classification by learning tasks in parallel. Motivated by the success of multi-task learning [Caruana, 1997], there are several neural network based NLP models [Collobert and Weston, 2008; Liu *et al.*, 2015b] utilize multitask learning to jointly learn several tasks with the aim of mutual benefit. The basic multi-task architectures of these models are to share some lower layers to determine common features. After the shared layers, the remaining layers are split into the multiple specific tasks.

In this paper, we propose three different models of sharing information with recurrent neural network (RNN). All the related tasks are integrated into a single system which is trained jointly. The first model uses just one shared layer for all the tasks. The second model uses different layers for different tasks, but each layer can read information from other layers. The third model not only assigns one specific layer for each task, but also builds a shared layer for all the tasks. Besides, we introduce a gating mechanism to enable the model to selectively utilize the shared information. The entire network is trained jointly on all these tasks.

Experimental results on four text classification tasks show that the joint learning of multiple related tasks together can improve the performance of each task relative to learning them separately.

Our contributions are of two-folds:

- First, we propose three multi-task architectures for RNN. Although the idea of multi-task learning is not new, our work is novel to integrate RNN into the multi-learning framework, which learns to map arbitrary text into semantic vector representations with both task-specific and shared layers.

- Second, we demonstrate strong results on several text classification tasks. Our multi-task models outperform most of state-of-the-art baselines.

---

*Corresponding author.

## 2 Recurrent Neural Network for Specific-Task Text Classification

The primary role of the neural models is to represent the variable-length text as a fixed-length vector. These models generally consist of a projection layer that maps words, sub-word units or n-grams to vector representations (often trained beforehand with unsupervised methods), and then combine them with the different architectures of neural networks.

There are several kinds of models to model text, such as Neural Bag-of-Words (NBOW) model, recurrent neural network (RNN) [Chung *et al.*, 2014], recursive neural network (RecNN) [Socher *et al.*, 2012; Socher *et al.*, 2013] and convolutional neural network (CNN) [Collobert *et al.*, 2011; Kalchbrenner *et al.*, 2014]. These models take as input the embeddings of words in the text sequence, and summarize its meaning with a fixed length vectorial representation.

Among them, recurrent neural networks (RNN) are one of the most popular architectures used in NLP problems because their recurrent structure is very suitable to process the variable-length text.

### 2.1 Recurrent Neural Network

A recurrent neural network (RNN) [Elman, 1990] is able to process a sequence of arbitrary length by recursively applying a transition function to its *internal hidden state vector* $h_t$ of the input sequence. The activation of the hidden state $h_t$ at time-step $t$ is computed as a function $f$ of the current input symbol $\mathbf{x}_t$ and the previous hidden state $\mathbf{h}_{t-1}$

$$\mathbf{h}_t = \begin{cases} 0 & t = 0 \\ f(\mathbf{h}_{t-1}, \mathbf{x}_t) & \text{otherwise} \end{cases} \quad (1)$$

It is common to use the state-to-state transition function $f$ as the composition of an element-wise nonlinearity with an affine transformation of both $\mathbf{x}_t$ and $\mathbf{h}_{t-1}$.

Traditionally, a simple strategy for modeling sequence is to map the input sequence to a fixed-sized vector using one RNN, and then to feed the vector to a softmax layer for classification or other tasks [Cho *et al.*, 2014].

Unfortunately, a problem with RNNs with transition functions of this form is that during training, components of the gradient vector can grow or decay exponentially over long sequences [Hochreiter *et al.*, 2001; Hochreiter and Schmidhuber, 1997]. This problem with *exploding* or *vanishing gradients* makes it difficult for the RNN model to learn long-distance correlations in a sequence.

Long short-term memory network (LSTM) was proposed by [Hochreiter and Schmidhuber, 1997] to specifically address this issue of learning long-term dependencies. The LSTM maintains a separate memory cell inside it that updates and exposes its content only when deemed necessary. A number of minor modifications to the standard LSTM unit have been made. While there are numerous LSTM variants, here we describe the implementation used by Graves [2013].

We define the LSTM *units* at each time step $t$ to be a collection of vectors in $\mathbb{R}^d$: an *input gate* $\mathbf{i}_t$, a *forget gate* $\mathbf{f}_t$, an *output gate* $\mathbf{o}_t$, a *memory cell* $\mathbf{c}_t$ and a hidden state $\mathbf{h}_t$. $d$ is the number of the LSTM units. The entries of the gating vectors $\mathbf{i}_t$, $\mathbf{f}_t$ and $\mathbf{o}_t$ are in $[0, 1]$. The LSTM transition equations
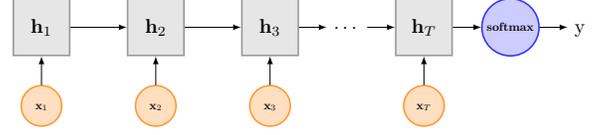


Figure 1: Recurrent Neural Network for Classification

are the following:

$$\mathbf{i}_t = \sigma(\mathbf{W}_i\mathbf{x}_t + \mathbf{U}_i\mathbf{h}_{t-1} + \mathbf{V}_i\mathbf{c}_{t-1}), \quad (2)$$
$$\mathbf{f}_t = \sigma(\mathbf{W}_f\mathbf{x}_t + \mathbf{U}_f\mathbf{h}_{t-1} + \mathbf{V}_f\mathbf{c}_{t-1}), \quad (3)$$
$$\mathbf{o}_t = \sigma(\mathbf{W}_o\mathbf{x}_t + \mathbf{U}_o\mathbf{h}_{t-1} + \mathbf{V}_o\mathbf{c}_t), \quad (4)$$
$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_c\mathbf{x}_t + \mathbf{U}_c\mathbf{h}_{t-1}), \quad (5)$$
$$\mathbf{c}_t = \mathbf{f}_t^i \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t, \quad (6)$$
$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t), \quad (7)$$

where $x_t$ is the input at the current time step, $\sigma$ denotes the logistic sigmoid function and $\odot$ denotes elementwise multiplication. Intuitively, the forget gate controls the amount of which each unit of the memory cell is erased, the input gate controls how much each unit is updated, and the output gate controls the exposure of the internal memory state.

### 2.2 Task-Specific Output Layer

In a single specific task, a simple strategy is to map the input sequence to a fixed-sized vector using one RNN, and then to feed the vector to a softmax layer for classification or other tasks.

Given a text sequence $x = \{x_1, x_2, \cdots, x_T\}$, we first use a lookup layer to get the vector representation (embeddings) $\mathbf{x}_i$ of the each word $x_i$. The output at the last moment $\mathbf{h}_T$ can be regarded as the representation of the whole sequence, which has a fully connected layer followed by a softmax non-linear layer that predicts the probability distribution over classes.

Figure 1 shows the unfolded RNN structure for text classification.

The parameters of the network are trained to minimise the cross-entropy of the predicted and true distributions.

$$L(\hat{\mathbf{y}}, \mathbf{y}) = -\sum_{i=1}^{N}\sum_{j=1}^{C} \mathbf{y}_i^j \log(\hat{\mathbf{y}}_i^j), \quad (8)$$

where $\mathbf{y}_i^j$ is the ground-truth label; $\hat{y}_i^j$ is prediction probabilities; $N$ denotes the number of training samples and $C$ is the class number.

## 3 Three Sharing Models for RNN based Multi-Task Learning

Most existing neural network methods are based on supervised training objectives on a single task [Collobert *et al.*, 2011; Socher *et al.*, 2013; Kalchbrenner *et al.*, 2014]. These methods often suffer from the limited amounts of training data. To deal with this problem, these models often involve an unsupervised pre-training phase. This unsupervised pre-training is effective to improve the final performance, but it does not directly optimize the desired task.

(a) Model-I: Uniform-Layer Architecture



(b) Model-II: Coupled-Layer Architecture
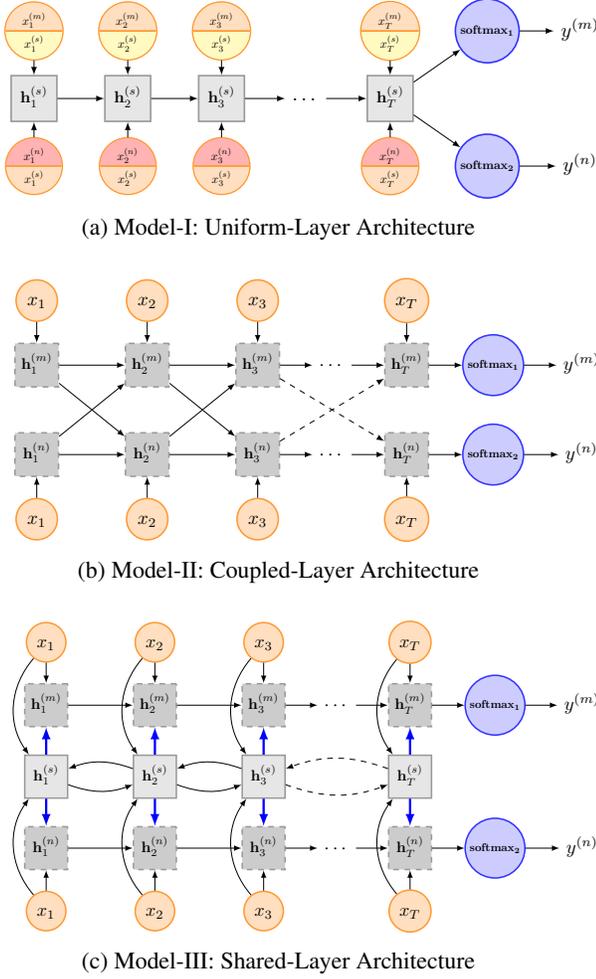


(c) Model-III: Shared-Layer Architecture

Figure 2: Three architectures for modelling text with multi-task learning.

Motivated by the success of multi-task learning [Caruana, 1997], we propose three multi-task models to leverage supervised data from many related tasks. Deep neural model is well suited for multi-task learning since the features learned from a task may be useful for other tasks. Figure 2 gives an illustration of our proposed models.

**Model-I: Uniform-Layer Architecture**   In Model-I, the different tasks share a same LSTM layer and an embedding layer besides their own embedding layers.

For task $m$, the input $\hat{\mathbf{x}}_t$ consists of two parts:

$$\hat{\mathbf{x}}_t^{(m)} = \mathbf{x}_t^{(m)} \oplus \mathbf{x}_t^{(s)}, \qquad (9)$$

where $\mathbf{x}_t^{(m)}$, $\mathbf{x}_t^{(s)}$ denote the task-specific and shared word embeddings respectively, $\oplus$ denotes the concatenation operation.

The LSTM layer is shared for all tasks. The final sequence representation for task $m$ is the output of LSMT at step $T$.

$$\mathbf{h}_T^{(m)} = LSTM(\hat{\mathbf{x}}^{(m)}). \qquad (10)$$

**Model-II: Coupled-Layer Architecture**   In Model-II, we assign a LSTM layer for each task, which can use the information for the LSTM layer of the other task.

Given a pair of tasks $(m, n)$, each task has own LSTM in the task-specific model. We denote the outputs at step $t$ of two coupled LSTM layer are $\mathbf{h}_t^{(m)}$ and $\mathbf{h}_t^{(n)}$.

To better control signals flowing from one task to another task, we use a global gating unit which endows the model with the capability of deciding how much information it should accept. We re-define Eqs. (5) and the new memory content of an LSTM at $m$-th task is computed by:

$$\tilde{\mathbf{c}}_t^{(m)} = \tanh\left(\mathbf{W}_c^{(m)}\mathbf{x}_t + \sum_{i \in \{m,n\}} \mathbf{g}^{(i \to m)} U_c^{(i \to m)} \mathbf{h}_{t-1}^{(i)}\right)$$
(11)

where $\mathbf{g}^{(i \to m)} = \sigma(\mathbf{W}_g^{(m)}\mathbf{x}_t + \mathbf{U}_g^{(i)}\mathbf{h}_{t-1}^{(i)})$. The other settings are same to the standard LSTM.

This model can be used to jointly learning for every two tasks. We can get two task specific representations $\mathbf{h}_T^{(m)}$ and $\mathbf{h}_T^{(n)}$ for tasks $m$ and $n$ receptively.

**Model-III: Shared-Layer Architecture**   Model-III also assigns a separate LSTM layer for each task, but introduces a bidirectional LSTM layer to capture the shared information for all the tasks.

We denote the outputs of the forward and backward LSTMs at step $t$ as $\overrightarrow{\mathbf{h}}_t^{(s)}$ and $\overleftarrow{\mathbf{h}}_t^{(s)}$ respectively. The output of shared layer is $\mathbf{h}_t^{(s)} = \overrightarrow{\mathbf{h}}_t^{(s)} \oplus \overleftarrow{\mathbf{h}}_t^{(s)}$.

To enhance the interaction between task-specific layers and the shared layer, we use gating mechanism to endow the neurons in task-specific layer with the ability to accept or refuse the information passed by the neuron in shared layers. Unlike Model-II, we compute the new state for LSTM as follows:

$$\tilde{\mathbf{c}}_t^{(m)} = \tanh\left(\mathbf{W}_c^{(m)}\mathbf{x}_t + \mathbf{g}^{(m)}\mathbf{U}_c^{(m)}\mathbf{h}_{t-1}^{(m)} + \mathbf{g}^{(s \to m)}\mathbf{U}_c^{(s)}\mathbf{h}_t^{(s)}\right),$$
(12)

where $\mathbf{g}^{(m)} = \sigma(\mathbf{W}_g^{(m)}\mathbf{x}_t + \mathbf{U}_g^{(m)}\mathbf{h}_{t-1}^{(m)})$ and $\mathbf{g}^{(s \to m)} = \sigma(\mathbf{W}_g^{(m)}\mathbf{x}_t + \mathbf{U}_g^{(s \to m)}\mathbf{h}_t^{(s)})$.

## 4   Training

The task-specific representations, which emittd by the muti-task architectures of all of the above, are ultimately fed into different output layers, which are also task-specific.

$$\hat{\mathbf{y}}^{(m)} = \text{softmax}(\mathbf{W}^{(m)}\mathbf{h}^{(m)} + \mathbf{b}^{(m)}), \qquad (13)$$

where $\hat{\mathbf{y}}^{(m)}$ is prediction probabilities for task $m$, $\mathbf{W}^{(m)}$ is the weight which needs to be learned, and $\mathbf{b}^{(m)}$ is a bias term.

Our global cost function is the linear combination of cost function for all joints.

$$\phi = \sum_{m=1}^{M} \lambda_m L(\hat{y}^{(m)}, y^{(m)}) \qquad (14)$$

where $\lambda_m$ is the weights for each task $m$ respectively.

| Dataset | Type | Train Size | Dev. Size | Test Size | Class | Averaged Length | Vocabulary Size |
|---------|------|-----------|-----------|-----------|-------|-----------------|-----------------|
| SST-1 | Sentence | 8544 | 1101 | 2210 | 5 | 19 | 18K |
| SST-2 | Sentence | 6920 | 872 | 1821 | 2 | 18 | 15K |
| SUBJ | Sentence | 9000 | - | 1000 | 2 | 21 | 21K |
| IMDB | Document | 25,000 | - | 25,000 | 2 | 294 | 392K |

Table 1: Statistics of the four datasets used in this paper.

It is worth noticing that labeled data for training each task can come from completely different datasets. Following [Collobert and Weston, 2008], the training is achieved in a stochastic manner by looping over the tasks:

1. Select a random task.
2. Select a random training example from this task.
3. Update the parameters for this task by taking a gradient step with respect to this example.
4. Go to 1.

**Fine Tuning**  For model-I and model-III, there is a shared layer for all the tasks. Thus, after the joint learning phase, we can use a fine tuning strategy to further optimize the performance for each task.

**Pre-training of the shared layer with neural language model**  For model-III, the shared layer can be initialized by an unsupervised pre-training phase. Here, for the shared LSTM layer in Model-III, we initialize it by a language model [Bengio *et al.*, 2007], which is trained on all the four task dataset.

## 5  Experiment

In this section, we investigate the empirical performances of our proposed three models on four related text classification tasks and then compare it to other state-of-the-art models.

### 5.1  Datasets

To show the effectiveness of multi-task learning, we choose four different text classification tasks about movie review. Each task have own dataset, which is briefly described as follows.

- **SST-1** The movie reviews with five classes (negative, somewhat negative, neutral, somewhat positive, positive) in the Stanford Sentiment Treebank[1] [Socher *et al.*, 2013].

- **SST-2** The movie reviews with binary classes. It is also from the Stanford Sentiment Treebank.

- **SUBJ** Subjectivity data set where the goal is to classify each instance (snippet) as being subjective or objective. [Pang and Lee, 2004]

- **IMDB** The IMDB dataset[2] consists of 100,000 movie reviews with binary classes [Maas *et al.*, 2011]. One key aspect of this dataset is that each movie review has several sentences.

---

[1] http://nlp.stanford.edu/sentiment.
[2] http://ai.stanford.edu/~amaas/data/sentiment/

The first three datasets are sentence-level, and the last dataset is document-level. The detailed statistics about the four datasets are listed in Table 1.

### 5.2  Hyperparameters and Training

The network is trained with backpropagation and the gradient-based optimization is performed using the Adagrad update rule [Duchi *et al.*, 2011]. In all of our experiments, the word embeddings are trained using word2vec [Mikolov *et al.*, 2013] on the Wikipedia corpus (1B words). The vocabulary size is about 500,000. The word embeddings are fine-tuned during training to improve the performance [Collobert *et al.*, 2011]. The other parameters are initialized by randomly sampling from uniform distribution in [-0.1, 0.1]. The hyperparameters which achieve the best performance on the development set will be chosen for the final evaluation. For datasets without development set, we use 10-fold cross-validation (CV) instead.

The final hyper-parameters are as follows. The embedding size for specific task and shared layer are 64. For Model-I, there are two embeddings for each word, and both their sizes are 64. The hidden layer size of LSTM is 50. The initial learning rate is 0.1. The regularization weight of the parameters is $10^{-5}$.

### 5.3  Effect of Multi-task Training

| Model | SST-1 | SST-2 | SUBJ | IMDB | Avg$\Delta$ |
|-------|-------|-------|------|------|-------------|
| Single Task | 45.9 | 85.8 | 91.6 | 88.5 | - |
| Joint Learning | 46.5 | 86.7 | 92.0 | 89.9 | +0.8 |
| + Fine Tuning | **48.5** | **87.1** | **93.4** | **90.8** | +2.0 |

Table 2: Results of the uniform-layer architecture.

| Model | SST-1 | SST-2 | SUBJ | IMDB | Avg$\Delta$ |
|-------|-------|-------|------|------|-------------|
| Single Task | 45.9 | 85.8 | 91.6 | 88.5 | - |
| SST1-SST2 | **48.9** | **87.4** | - | - | +2.3 |
| SST1-SUBJ | 46.3 | - | 92.2 | - | +0.5 |
| SST1-IMDB | 46.9 | - | - | 89.5 | +1.0 |
| SST2-SUBJ | - | 86.5 | 92.5 | - | +0.8 |
| SST2-IMDB | - | 86.8 | - | **89.8** | +1.2 |
| SUBJ-IMDB | - | - | **92.7** | 89.3 | +0.9 |

Table 3: Results of the coupled-layer architecture.

We first compare our our proposed models with the standard LSTM for single task classification. We use the implementation of Graves [2013]. The unfolded illustration is shown in Figure 1.

| Model | SST-1 | SST-2 | SUBJ | IMDB | Avg$\Delta$ |
|---|---|---|---|---|---|
| Single Task | 45.9 | 85.8 | 91.6 | 88.5 | - |
| Joint Learning | 47.1 | 87.0 | 92.5 | 90.7 | +1.4 |
| + LM | 47.9 | 86.8 | 93.6 | 91.0 | +1.9 |
| + Fine Tuning | **49.6** | **87.9** | **94.1** | **91.3** | +2.8 |

Table 4: Results of the shared-layer architecture.

| Model | SST-1 | SST-2 | SUBJ | IMDB |
|---|---|---|---|---|
| NBOW | 42.4 | 80.5 | 91.3 | 83.62 |
| MV-RNN | 44.4 | 82.9 | - | - |
| RNTN | 45.7 | 85.4 | - | - |
| DCNN | 48.5 | 86.8 | - | - |
| PV | 44.6 | 82.7 | 90.5 | **91.7** |
| Tree-LSTM | **50.6** | 86.9 | - | - |
| Multi-Task | 49.6 | **87.9** | **94.1** | 91.3 |

Table 5: Results of shared-layer multi-task model against state-of-the-art neural models.

Table 2-4 show the classification accuracies on the four datasets. The second line ("Single Task") of each table shows the result of the standard LSTM for each individual task.

**Uniform-layer Architecture** For the first uniform-layer architecture, we train the model on four datasets simultaneously. The LSTM layer is shared across all the tasks. The average improvement of the performances on four datasets is 0.8%. With the further fine-tuning phase, the improvement achieves 2.0% on average.

**Coupled-layer Architecture** For the second coupled-layer architecture, the information is shared with a pair of tasks. Therefore, there are six combinations for the four datasets. We train six models on the different pairs of datasets. We can find that the pair-wise joint learning also improves the performances. The more relevant the tasks are, the more significant the improvements are. Since SST-1 and SST-2 are from the same corpus, their improvements are more significant than the other combinations. The improvement is 2.3% on average with simultaneously learning on SST-1 and SST-2.

**Shared-layer Architecture** The shared-layer architecture is more general than uniform-layer architecture. Besides a shared layer for all the tasks, each task has own task-specific layer. As shown in Table 4, we can see that the average improvement of the performances on four datasets is 1.4%, which is better than the uniform-layer architecture. We also investigate the strategy of unsupervised pre-training towards shared LSTM layer. With the LM pre-training, the performance is improved by an extra 0.5% on average. Besides, the further fine-tuning can significantly improve the performance by an another 0.9%.

To recap, all our proposed models outperform the baseline of single-task learning. The shared-layer architecture gives the best performances. Moreover, compared with vanilla LSTM, our proposed three models don't cause much extra computational cost while converge faster. In our experiment, the most complicated model-III, costs 2.5 times as long as vanilla LSTM.

### 5.4 Comparisons with State-of-the-art Neural Models

We compare our model with the following models:

- **NBOW** The NBOW sums the word vectors and applies a non-linearity followed by a softmax classification layer.

- **MV-RNN** Matrix-Vector Recursive Neural Network with parse trees [Socher *et al.*, 2012].

- **RNTN** Recursive Neural Tensor Network with tensor-based feature function and parse trees [Socher *et al.*, 2013].

- **DCNN** Dynamic Convolutional Neural Network with dynamic k-max pooling [Kalchbrenner *et al.*, 2014].

- **PV** Logistic regression on top of paragraph vectors [Le and Mikolov, 2014]. Here, we use the popular open source implementation of PV in Gensim[3].

- **Tree-LSTM** A generalization of LSTMs to tree-structured network topologies. [Tai *et al.*, 2015]

Table 5 shows the performance of the shared-layer architecture compared with the competitor models, which shows our model is competitive for the neural-based state-of-the-art models.

Although Tree-LSTM outperforms our model on SST-1, it needs an external parser to get the sentence topological structure. It is worth noticing that our models are compatible with the other RNN based models. For example, we can easily extend our models to incorporate the Tree-LSTM model.

### 5.5 Case Study

To get an intuitive understanding of what is happening when we use the single LSTM or the shared-layer LSTM to predict the class of text, we design an experiment to analyze the output of the single LSTM and the shared-layer LSTM at each time step. We sample two sentences from the SST-2 test dataset, and the changes of the predicted sentiment score at different time steps are shown in Figure 3. To get more insights into how the shared structures influences the specific task. We observe the activation of global gates $g^{(s)}$, which controls signals flowing from one shared LSTM layer to task-spcific layer, to understand the behaviour of neurons. We plot evolving activation of global gates $g^{(s)}$ through time and sort the neurons according to their activations at the last time step.

For the sentence "`A merry movie about merry period people's life.`", which has a positive sentiment, while the standard LSTM gives a wrong prediction. The reason can be inferred from the activation of global gates $g^{(s)}$. As shown in Figure 3-(c), we can see clearly the neurons are activated much when they take input as "`merry`", which indicates the task-specific layer takes much information from shared layer towards the word "`merry`", and this ultimately makes the model give a correct prediction.

---

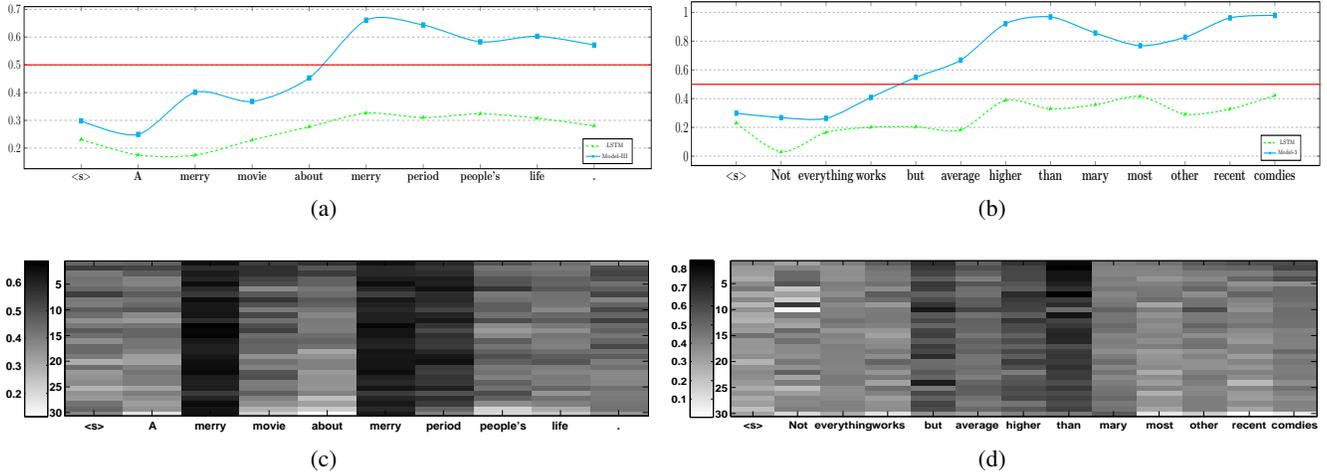[3]`https://github.com/piskvorky/gensim/`

Figure 3: (a)(b) The change of the predicted sentiment score at different time steps. Y-axis represents the sentiment score, while X-axis represents the input words in chronological order. The red horizontal line gives a border between the positive and negative sentiments. (c)(d) Visualization of the global gate's ($g^{(s)}$) activation.

Another case "`Not everything works, but the average is higher than in Mary and most other recent comedies.`" is positive and has a little complicated semantic composition. As shown in Figure 3-(b,d), simple LSTM cannot capture the structure of "`but ... higher than`" while our model is sensitive to it, which indicates the shared layer can not only enrich the meaning of certain words, but can teach some information of structure to specific task.

## 5.6 Error Analysis

We analyze the bad cases induced by our proposed shared-layer model on SST-2 dataset. Most of the bad cases can be generalized into two categories

**Complicated Sentence Structure** Some sentences involved complicated structure can not be handled properly, such as double negation "`it never fails to engage us.`" and subjunctive sentences "`Still, I thought it could have been more.`". To solve these cases, some architectural improvements are necessary, such as tree-based LSTM [Tai *et al.*, 2015].

**Sentences Required Reasoning** The sentiments of some sentences can be mislead if only considering the literal meaning. For example, the sentence "`I tried to read the time on my watch.`" expresses negative attitude towards a movie, which can be understood correctly by reasoning based on common sense.

## 6 Related Work

Neural networks based multi-task learning has proven effective in many NLP problems [Collobert and Weston, 2008; Liu *et al.*, 2015b].

Collobert and Weston [2008] used a shared representation for input words and solve different traditional NLP tasks such as part-of-Speech tagging and semantic role labeling within one framework. However, only one lookup table is shared, and the other lookup-tables and layers are task specific. To deal with the variable-length text sequence, they used window-based method to fix the input size.

Liu *et al.* [2015b] developed a multi-task DNN for learning representations across multiple tasks. Their multi-task DNN approach combines tasks of query classification and ranking for web search. But the input of the model is bag-of-word representation, which lose the information of word order.

Different with the two above methods, our models are based on recurrent neural network, which is better to model the variable-length text sequence.

More recently, several multi-task encoder-decoder networks were also proposed for neural machine translation [Dong *et al.*, 2015; Firat *et al.*, 2016], which can make use of cross-lingual information. Unlike these works, in this paper we design three architectures, which can control the information flow between shared layer and task-specific layer flexibly, thus obtaining better sentence representations.

## 7 Conclusion and Future Work

In this paper, we introduce three RNN based architectures to model text sequence with multi-task learning. The differences among them are the mechanisms of sharing information among the several tasks. Experimental results show that our models can improve the performances of a group of related tasks by exploring common features.

In future work, we would like to investigate the other sharing mechanisms of the different tasks.

# References

[Bengio *et al.*, 2007] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, et al. Greedy layer-wise training of deep networks. *Advances in neural information processing systems*, 19:153, 2007.

[Caruana, 1997] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.

[Cho *et al.*, 2014] Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of EMNLP*, 2014.

[Chung *et al.*, 2014] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.

[Collobert and Weston, 2008] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of ICML*, 2008.

[Collobert *et al.*, 2011] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537, 2011.

[Dong *et al.*, 2015] Daxiang Dong, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. Multi-task learning for multiple language translation. In *Proceedings of the ACL*, 2015.

[Duchi *et al.*, 2011] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159, 2011.

[Elman, 1990] Jeffrey L Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.

[Firat *et al.*, 2016] Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. Multi-way, multilingual neural machine translation with a shared attention mechanism. *arXiv preprint arXiv:1601.01073*, 2016.

[Graves, 2013] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.

[Hochreiter and Schmidhuber, 1997] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[Hochreiter *et al.*, 2001] Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, and Jürgen Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.

[Kalchbrenner *et al.*, 2014] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modelling sentences. In *Proceedings of ACL*, 2014.

[Le and Mikolov, 2014] Quoc V. Le and Tomas Mikolov. Distributed representations of sentences and documents. In *Proceedings of ICML*, 2014.

[Liu *et al.*, 2015a] PengFei Liu, Xipeng Qiu, Xinchi Chen, Shiyu Wu, and Xuanjing Huang. Multi-timescale long short-term memory neural network for modelling sentences and documents. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2015.

[Liu *et al.*, 2015b] Xiaodong Liu, Jianfeng Gao, Xiaodong He, Li Deng, Kevin Duh, and Ye-Yi Wang. Representation learning using multi-task deep neural networks for semantic classification and information retrieval. In *NAACL*, 2015.

[Maas *et al.*, 2011] Andrew L Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the ACL*, pages 142–150, 2011.

[Mikolov *et al.*, 2013] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

[Pang and Lee, 2004] Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of ACL*, 2004.

[Socher *et al.*, 2011] Richard Socher, Jeffrey Pennington, Eric H Huang, Andrew Y Ng, and Christopher D Manning. Semi-supervised recursive autoencoders for predicting sentiment distributions. In *Proceedings of EMNLP*, 2011.

[Socher *et al.*, 2012] Richard Socher, Brody Huval, Christopher D Manning, and Andrew Y Ng. Semantic compositionality through recursive matrix-vector spaces. In *Proceedings of EMNLP*, pages 1201–1211, 2012.

[Socher *et al.*, 2013] Richard Socher, Alex Perelygin, Jean Y Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of EMNLP*, 2013.

[Tai *et al.*, 2015] Kai Sheng Tai, Richard Socher, and Christopher D Manning. Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint arXiv:1503.00075*, 2015.

[Turian *et al.*, 2010] Joseph Turian, Lev Ratinov, and Yoshua Bengio. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of ACL*, 2010.