

# SGM: Sequence Generation Model for Multi-Label Classification

Pengcheng Yang (speaker), Xu Sun, Wei Li  
Shuming Ma, Wei Wu, Houfeng Wang

**Peking University, Beijing, China**

24<sup>th</sup> Aug 2018

# Outline

- ① Introduction to Multi-Label Classification
- ② Proposal: Sequence Generation Model
- ③ Experiments and Analysis
- ④ Conclusion

# Introduction to Multi-Label Classification

# What is Multi-Label Classification?

## ① Definition:

- Assign multiple labels to each sample in the dataset.

# What is Multi-Label Classification?

## 1 Definition:

- Assign multiple labels to each sample in the dataset.

## 2 Example:

[6] [arXiv:1807.07545](#) [pdf, other]

### Rearranging the Familiar: Testing Compositional Generalization in Recurrent Networks

João Loula, Marco Baroni, Brenden M. Lake

Subjects: [Computation and Language \(cs.CL\)](#); [Artificial Intelligence \(cs.AI\)](#); [Machine Learning \(cs.LG\)](#)

[7] [arXiv:1807.07520](#) [pdf, ps, other]

### Statistical Model Compression for Small-Footprint Natural Language Understanding

Grant P. Strimel, Kanthashree Mysore Sathyendra, Stanislav Peshterliev

Comments: Interspeech 2018

Subjects: [Computation and Language \(cs.CL\)](#)

[8] [arXiv:1807.07517](#) [pdf, other]

### Using Deep Neural Networks to Translate Multi-lingual Threat Intelligence

Priyanka Ranade, Sudip Mittal, Anupam Joshi, Karuna Joshi

Subjects: [Computation and Language \(cs.CL\)](#); [Cryptography and Security \(cs.CR\)](#)

# What is Multi-Label Classification?

## 1 Definition:

- Assign multiple labels to each sample in the dataset.

## 2 Example:

[6] [arXiv:1807.07545](#) [pdf, other]

### Rearranging the Familiar: Testing Compositional Generalization in Recurrent Networks

João Loula, Marco Baroni, Brenden M. Lake

Subjects: [Computation and Language \(cs.CL\)](#); [Artificial Intelligence \(cs.AI\)](#); [Machine Learning \(cs.LG\)](#)

[7] [arXiv:1807.07520](#) [pdf, ps, other]

### Statistical Model Compression for Small-Footprint Natural Language Understanding

Grant P. Strimel, Kanthashree Mysore Sathyendra, Stanislav Peshterliev

Comments: Interspeech 2018

Subjects: [Computation and Language \(cs.CL\)](#)

[8] [arXiv:1807.07517](#) [pdf, other]

### Using Deep Neural Networks to Translate Multi-lingual Threat Intelligence

Priyanka Ranade, Sudip Mittal, Anupam Joshi, Karuna Joshi

Subjects: [Computation and Language \(cs.CL\)](#); [Cryptography and Security \(cs.CR\)](#)

## 3 Applications:

- Text categorization, information retrieval, and so on.

# Background

## Previous work:

- ① Can't capture **label correlations** very well or is **computationally intractable**.
  - **Label correlations:** Some labels are closely correlated.

# Background

## Previous work:

- ① Can't capture **label correlations** very well or is **computationally intractable**.
  - **Label correlations:** Some labels are closely correlated.
- ② Ignore **differences in the contributions** of textual content when predicting different labels.

|  |
|--|
| • Generating descriptions for <b>videos</b> has many applications including human <b>robot</b> interaction.                                |
| • Many methods for <b>image captioning</b> rely on pre-trained <b>object classifier CNN</b> and Long Short Term Memory recurrent networks. |
| • How to learn <b>robust visual classifiers</b> from the weak annotations of the sentence descriptions.                                    |

(a) Visual analysis when the SGM model predicts "CV".

|  |
|--|
| • Generating descriptions for <b>videos</b> has many applications including human robot interaction.                                       |
| • Many methods for image captioning rely on pre-trained object classifier CNN and <b>Long</b> Short Term <b>Memory recurrent</b> networks. |
| • How to learn <b>robust visual classifiers</b> from the weak <b>annotations</b> of the <b>sentence</b> descriptions.                      |

(b) Visual analysis when the SGM model predicts "CL".

Figure 1: Visualization of attention.



# Proposal: Sequence Generation Model

# Proposal: Sequence Generation Model

**Transform classification task into generation task.**

**① Key ideas:**

- View the **text** as the **source language** and the **label** as **target language**.
- Base on **sequence-to-sequence** model.

# Proposal: Sequence Generation Model

## Transform classification task into generation task.

### ① Key ideas:

- View the **text** as the **source language** and the **label** as **target language**.
- Base on **sequence-to-sequence** model.

### ② Advantages:

- **Capture label correlations:** Generate labels **sequentially**, and predict the next label based on its previously generated labels.
- **Consider differences in contributions of textual content:** Apply the **attention** mechanism.

# Proposal: Sequence Generation Model

## Difficulties and solutions:

### ① Repeated labels:

- Use the **masked softmax** layer to **smooth** the probability distribution.

# Proposal: Sequence Generation Model

## Difficulties and solutions:

### ① Repeated labels:

- Use the **masked softmax** layer to **smooth** the probability distribution.

### ② Exposure bias:

- Use the **adaptive gate** to introduce the **global information** of previous time-steps.

# Proposal: Sequence Generation Model

## Difficulties and solutions:

### ① Repeated labels:

- Use the **masked softmax** layer to **smooth** the probability distribution.

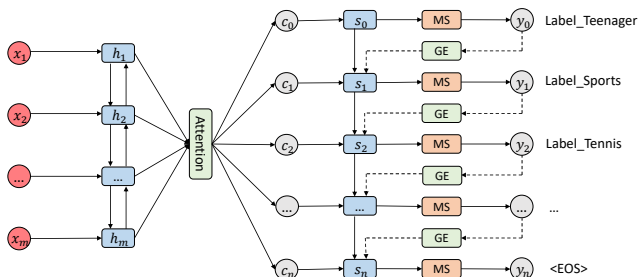
### ② Exposure bias:

- Use the **adaptive gate** to introduce the **global information** of previous time-steps.

### ③ Sequence order:

- **Sort** the label sequence of each sample according to the **frequency** of labels and **high-frequency** labels are placed in the **front**.

# Proposal: Sequence Generation Model



**Figure 2:** The overview of SGM with global embedding. MS denotes the masked softmax layer. GE denotes the global embedding.

- The proposed model is based on the **Seq2Seq** model, which consists of an **encoder** and a **decoder with global embedding**.

# Proposal: Masked Softmax Layer

**Masked softmax layer:** Prevent the decoder from predicting **repeated labels**.

- 1  $y_t$  is the probability distribution over the label space  $\mathcal{L}$  at time-step  $t$ .

$$y_t = \text{softmax}(o_t + l_t) \quad (1)$$

- 2  $l_t \in \mathbb{R}^L$  is the mask vector.

$$(l_t)_i = \begin{cases} -\infty & \text{if the label } l_i \text{ has been predicted.} \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$



# Proposal: Global Embedding

**Global embedding:** Introduce the **global information** of previous time-steps to **alleviate the exposure bias**.

$$\bar{e} = \sum_{i=1}^L y_{t-1}^{(i)} e_i \quad (3)$$

$$g(y_t) = (1 - H) \odot e + H \odot \bar{e} \quad (4)$$

$$H = \mathbf{W}_1 e + \mathbf{W}_2 \bar{e} \quad (5)$$

- ①  $e$  is the embedding vector of the label which has the highest probability under distribution  $y_{t-1}$ .
- ②  $e_i$  is the embedding vector of the  $i$ -th label.
- ③  $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{L \times L}$  are weight matrices.

# Experiments and Analysis

# Datasets and Evaluation Metrics

## ① Datasets:

- **RCV1-V2:** Reuters Corpus Volume I.
- **AAPD:** Arxiv Academic Paper Dataset.

## ② Evaluation metrics:

- **Hamming loss** and **micro- $F_1$**  score are our main evaluation metrics.
- **Micro-precision** and **micro-recall** are also reported to assist the analysis.

# Results

| Models  | HL(-)         | P(+)         | R(+)         | F1(+)        |
|---------|---------------|--------------|--------------|--------------|
| BR      | 0.0086        | 0.904        | 0.816        | 0.858        |
| CC      | 0.0087        | 0.887        | 0.828        | 0.857        |
| LP      | 0.0087        | 0.896        | 0.824        | 0.858        |
| CNN     | 0.0089        | <b>0.922</b> | 0.798        | 0.855        |
| CNN-RNN | 0.0085        | 0.889        | 0.825        | 0.856        |
| SGM     | 0.0081        | 0.887        | 0.850        | 0.869        |
| + GE    | <b>0.0075</b> | 0.897        | <b>0.860</b> | <b>0.878</b> |

(a) Performance on RCV1-V2 test set.

| Models  | HL(-)         | P(+)         | R(+)         | F1(+)        |
|---------|---------------|--------------|--------------|--------------|
| BR      | 0.0316        | 0.644        | 0.648        | 0.646        |
| CC      | 0.0306        | 0.657        | 0.651        | 0.654        |
| LP      | 0.0312        | 0.662        | 0.608        | 0.634        |
| CNN     | 0.0256        | <b>0.849</b> | 0.545        | 0.664        |
| CNN-RNN | 0.0278        | 0.718        | 0.618        | 0.664        |
| SGM     | 0.0251        | 0.746        | 0.659        | 0.699        |
| + GE    | <b>0.0245</b> | 0.748        | <b>0.675</b> | <b>0.710</b> |

(b) Performance on AAPD test set.

**Table 1:** Comparison between our methods and all baselines on two datasets. GE denotes the global embedding. HL, P, R, and F1 denote hamming loss, micro-precision, micro-recall, and micro- $F_1$ , respectively.

# Exploration of Global Embedding

- ① **Goal:** Explore how the performance of our model is affected by the proportion between two kinds of embeddings.
- ② **Settings:**
  - Adaptive gate:

$$g(y_t) = (1 - H) \odot e + H \odot \bar{e} \quad (6)$$

$$H = \mathbf{W}_1 e + \mathbf{W}_2 \bar{e} \quad (7)$$

- Coefficient averaging:

$$g(y_t) = (1 - \lambda) * e + \lambda * \bar{e} \quad (8)$$

# Exploration of Global Embedding

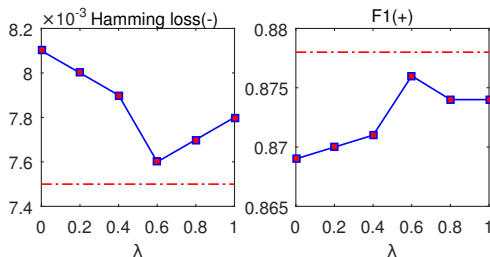


Figure 3: The performance of the SGM when using different  $\lambda$ . The red dotted line represents the results of using the adaptive gate.

# Exploration of Global Embedding

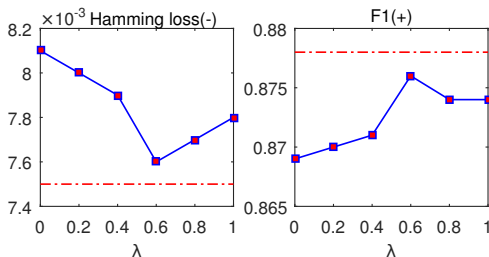


Figure 3: The performance of the SGM when using different  $\lambda$ . The red dotted line represents the results of using the adaptive gate.

- 1 The weighted average embedding contains **richer information**, leading to the improvement in the performance of the model.
- 2 The adaptive gate can automatically determine the **most appropriate  $\lambda$  value** according to the actual condition.

# Ablation Study

| Models             | HL(-)           | F1(+)          | Models             | HL(-)            | F1(+)          |
|--------------------|-----------------|----------------|--------------------|------------------|----------------|
| SGM                | 0.0081          | 0.869          | SGM + GE           | 0.0075           | 0.878          |
| <i>w/o mask</i>    | 0.0083(↓ 2.47%) | 0.866(↓ 0.35%) | <i>w/o mask</i>    | 0.0078(↓ 4.00%)  | 0.873(↓ 0.57%) |
| <i>w/o sorting</i> | 0.0084(↓ 3.70%) | 0.858(↓ 1.27%) | <i>w/o sorting</i> | 0.0083(↓ 10.67%) | 0.859(↓ 2.16%) |

(a) Ablation study for SGM.

(b) Ablation study for SGM with GE.

**Table 2:** Ablation study on the RCV1-V2 test set. GE denotes global embedding. ↓ indicates that the performance of the model is degraded.



# Ablation Study

| Models             | HL(-)           | F1(+)          | Models             | HL(-)            | F1(+)          |
|--------------------|-----------------|----------------|--------------------|------------------|----------------|
| SGM                | 0.0081          | 0.869          | SGM + GE           | 0.0075           | 0.878          |
| <i>w/o mask</i>    | 0.0083(↓ 2.47%) | 0.866(↓ 0.35%) | <i>w/o mask</i>    | 0.0078(↓ 4.00%)  | 0.873(↓ 0.57%) |
| <i>w/o sorting</i> | 0.0084(↓ 3.70%) | 0.858(↓ 1.27%) | <i>w/o sorting</i> | 0.0083(↓ 10.67%) | 0.859(↓ 2.16%) |

(a) Ablation study for SGM.

(b) Ablation study for SGM with GE.

**Table 2:** Ablation study on the RCV1-V2 test set. GE denotes global embedding. ↓ indicates that the performance of the model is degraded.

- 1 Sorting is important because humans need to **predefine the order** of output labels.
- 2 The mask module has little impact because **label cardinality is small**.

# Error Analysis

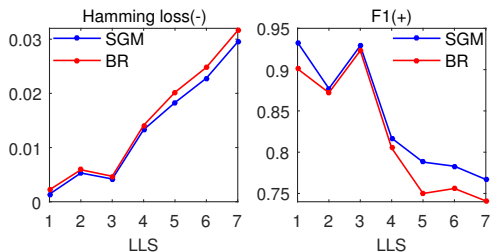


Figure 4: The performance of SGM on different subsets of the RCV1-V2 test set. LLS represents the length of label sequence of each sample in the subset.

# Error Analysis

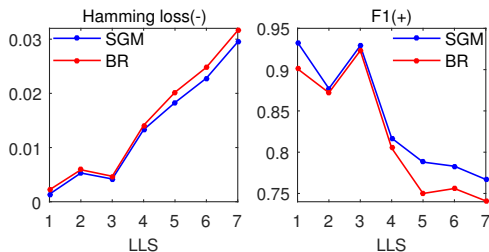


Figure 4: The performance of SGM on different subsets of the RCV1-V2 test set. LLS represents the length of label sequence of each sample in the subset.

- 1 The performance of all methods **deteriorates** when **LLS increases**.
- 2 The advantages of SGM are more **significant** when **LLS is large**.

# Visualization of Attention

|  |
|--|
| <ul style="list-style-type: none"> <li>Generating descriptions for videos has many applications including human robot interaction.</li> </ul>                                |
| <ul style="list-style-type: none"> <li>Many methods for image captioning rely on pre-trained object classifier CNN and Long Short Term Memory recurrent networks.</li> </ul> |
| <ul style="list-style-type: none"> <li>How to learn robust visual classifiers from the weak annotations of the sentence descriptions.</li> </ul>                             |

(a) Visual analysis when the SGM model predicts “CV”.

|  |
|--|
| <ul style="list-style-type: none"> <li>Generating descriptions for videos has many applications including human robot interaction.</li> </ul>                                |
| <ul style="list-style-type: none"> <li>Many methods for image captioning rely on pre-trained object classifier CNN and Long Short Term Memory recurrent networks.</li> </ul> |
| <ul style="list-style-type: none"> <li>How to learn robust visual classifiers from the weak annotations of the sentence descriptions.</li> </ul>                             |

(b) Visual analysis when the SGM model predicts “CL”.

**Figure 5:** An example abstract in the AAPD dataset, from which we extract three informative sentences. This abstract is assigned two labels: “CV” and “CL”. They denote computer vision and computational language, respectively.

# Visualization of Attention

|  |
|--|
| <ul style="list-style-type: none"> <li>Generating descriptions for videos has many applications including human robot interaction.</li> </ul>                                |
| <ul style="list-style-type: none"> <li>Many methods for image captioning rely on pre-trained object classifier CNN and Long Short Term Memory recurrent networks.</li> </ul> |
| <ul style="list-style-type: none"> <li>How to learn robust visual classifiers from the weak annotations of the sentence descriptions.</li> </ul>                             |

(a) Visual analysis when the SGM model predicts “CV”.

|  |
|--|
| <ul style="list-style-type: none"> <li>Generating descriptions for videos has many applications including human robot interaction.</li> </ul>                                |
| <ul style="list-style-type: none"> <li>Many methods for image captioning rely on pre-trained object classifier CNN and Long Short Term Memory recurrent networks.</li> </ul> |
| <ul style="list-style-type: none"> <li>How to learn robust visual classifiers from the weak annotations of the sentence descriptions.</li> </ul>                             |

(b) Visual analysis when the SGM model predicts “CL”.

**Figure 5:** An example abstract in the AAPD dataset, from which we extract three informative sentences. This abstract is assigned two labels: “CV” and “CL”. They denote computer vision and computational language, respectively.

- The attention mechanism can select the most informative words automatically when predicting different labels.

# Case Study

| Reference   | BR                    | SGM                                     | SGM + GE   |
|---|-----------------------|---|--|
| CCAT, <b>C15, C152</b> , C41, C411                      | CCAT, C15, C13        | CCAT, <b>C15, C152</b>                  | CCAT, <b>C15, C152</b> , C41, C411                   |
| CCAT, GCAT, ECAT, C31, GDIP, C13, C21, <b>E51, E512</b> | CCAT, GCAT, GDIP, E51 | CCAT, ECAT, GDIP, <b>E51, E512</b>      | CCAT, GCAT, ECAT, C31, GDIP, <b>E51, E512</b> , C312 |
| GCAT, ECAT, <b>G15, G154, G151, G155</b>                | GCAT, ECAT, GENV, G15 | GCAT, ECAT, E21, <b>G15, G154, G156</b> | GCAT, ECAT, E21, <b>G15, G154, G155</b>              |

**Figure 6:** Several examples of generated label sequences on the RCV1-V2 dataset. The red bold labels in each example indicate that they are highly correlated.

# Case Study

| Reference   | BR                    | SGM                                     | SGM + GE   |
|---|-----------------------|---|--|
| CCAT, <b>C15, C152</b> , C41, C411                      | CCAT, C15, C13        | CCAT, <b>C15, C152</b>                  | CCAT, <b>C15, C152</b> , C41, C411                   |
| CCAT, GCAT, ECAT, C31, GDIP, C13, C21, <b>E51, E512</b> | CCAT, GCAT, GDIP, E51 | CCAT, ECAT, GDIP, <b>E51, E512</b>      | CCAT, GCAT, ECAT, C31, GDIP, <b>E51, E512</b> , C312 |
| GCAT, ECAT, <b>G15, G154, G151, G155</b>                | GCAT, ECAT, GENV, G15 | GCAT, ECAT, E21, <b>G15, G154, G156</b> | GCAT, ECAT, E21, <b>G15, G154, G155</b>              |

**Figure 6:** Several examples of generated label sequences on the RCV1-V2 dataset. The red bold labels in each example indicate that they are highly correlated.

- 1 The proposed SGM can **capture the correlations** between labels.
- 2 The SGM with global embedding predicts labels **more accurately**.

# Conclusion



# Conclusion

- ① The sequence generation model is able to **capture the correlations between labels** well.
- ② The attention mechanism can **select the most informative words automatically** when predicting different labels.
- ③ The global embedding can **alleviate exposure bias** by introducing the **global information** of previous time-steps.

- If there is any question, please contact Pengcheng Yang (yang\_pc@pku.edu.cn)
- The code and datasets are available at <https://github.com/lancopku/SGM>

**Thank you!**