

Summary of Research

Zhe Xu

Shanghai Key Laboratory of Data Science

Fudan University

Email: zxu14@fudan.edu.cn

Supervisor:
Yanghua Xiao (FDU)

Outline

- Incorporating Complicated Rules into Deep Generative Models
- Relation Extraction: From the Perspective of Implicit Relation Reasoning
- Hierarchical Conceptual Labelling
- Entity Typing with DNN and Word Embedding

*Above all are research projects mainly. If you have any interests about my engineering projects please contact me. 😊

1. Incorporating Complicated Rules into Deep Generative Models

1. Incorporating Complicated Rules into Deep Generative Models

Problem

Results of Generative Tasks are far from perfect, for example:



Generated digits have fissures (3 and 6 in this image) and noise points. This image has been selected to show the problem.

*"...a marble already & saies foorth
priam foorth anger anger fearing
norway anger mutiny noyse brought
downe."*

Generated Shakespeare text has word repetition.

Above all problems are easy for human to find. But if we want the model to learn we need to tune the model carefully and provide much more data.

1. Incorporating Complicated Rules into Deep Generative Models

Idea

GAN (generative adversarial network) will offer a feedback to the Generator through a network named Discriminator. Through training the Discriminator to distinguish the generated and original samples, improve the Generator at the same time.

Inspired by this, we designed another Rule Discriminator which is **trained to judge whether the results (no matter generated or original) is consistent with some Rules (e.g. No word repetition in a sentence.)**. Then we concatenated the Rule Discriminator with the Generator to improve the generated results.

It is easy for us to judge whether the results are qualified with some rules by some 'If.....Else.....', we called it **Rule Examiner**. See more detail in next page.

Method

- We do not have enough positive and negative samples at first and **only using Rule Discriminator to train Generator is not enough**. So the **original loss function should be maintained**.
- Since **the Rule Examiner is a part of this framework** (it need to examine the generated results in every epoch), if it works slowly (e.g. We use flood-fill algorithm to judge if there exists noise points or fissures in Hand-Written digits generation task), the training processing will be slow.
- We **do not use Rule Discriminator at first** since we do not have enough samples and **at the early training stage the generated results are useless** for Rule Discriminator to learn.(e.g. A blank image)
- The **Rules should be consistent with original loss**(e.g. Most Shakespeare text doesn't have continuous word repetition like '...anger anger...'). But we have noticed that if we let the model overemphasize the Rules like no noise points in Hand-Written digits, the generator will generate something wrapped together even like a solid white block. So we **set a hyper-parameter to adjust the loss function of Generator**.

Example1 :Hand-Written Digits

Structure

We used our framework into a Hand-Written Digits task (based on AC-GAN). And the structure of our Rule Discriminator is the same as the original Discriminator but without the Number label classification output.

Rule

In order to eliminate the noise points and fissures, we set the rule:

There should be only one white connected block. We examined it by flood-fill algorithm.

Example1 :Hand-Written Digits

Result

From the images we can find that there is less fissures after using our framework. We can get that more directly in our statistic results. P_R is the percentage that the results can pass the Rule Examiner (We have used it as part of our framework). P_H is the percentage that human marked the result 'good'.



Without our framework



With our framework

Metric	MNIST test data	AC-GAN	Our model
P_R	96.2%	96.5%	99.7%
P_H	57.1%	58.1%	74.4%

Statistic results

Example2: Chinese Question

Structure

Our baseline model is VAE-GAN, and the structure of Rule Discriminator is the same as the original one without a Predicate output.

Rule

Rules for Chinese question generation

1 Sentences should end with ‘#’ (a special character).

2 The subject should appear only once in a sentence.

3 There are no continuously repeated characters in the sentence.

4 The length of sentences should be more than 4 characters.

(A Chinese question should not be too short.)

5 The number of low frequency words should less than half of sentence length.

*We replace the subject in sentence with ‘~’, and our training sentences all end with ‘#’.

Example2: Chinese Question

Result

The definition of P_R and P_H are same as the previous ones. And we show some cases with translation or annotation in the table.

With our framework, the quality of the generated text or image has been improved a lot.

Metric	Original model	Our model
P_H	51.0%	57.0%
P_R for all rules	63.0%	74.2%
P_R for rule 1	98.2%	99.0%
P_R for rule 2	92.8%	96.8%
P_R for rule 3	76.2%	83.8%
P_R for rule 4	100.0%	100.0%
P_R for rule 5	99.8%	100.0%

Predicate	Question(without R)	Question(with R)
原料	~是什么原料	你知道~需要什么原料吗
Raw material	What is the raw material of	Do you know what raw material does ~ need
编剧	谁问是的是主剧的	你知道~的编剧是谁吗
Screenwriter	(unreadable and without ~)	Do you know who the screenwriter of ~ is
别名	你知道~这的形是么吗	~还叫什么
Alias	Do you know ~*** (continuous repeated words)	What is ~ also called
地点	~发生在什么地方	~发生在什么地方
Location	Where does ~ happen	Where does ~ happen
作者	《知》我~的者演是谁吗	你知道~这本书的作者是谁吗
author	(unreadable)	Do you know who the author of this book ~ is

2. Relation Extraction: From the Perspective of Implicit Relation Reasoning

Problem

Relation extraction (RE) is a fundamental task for text analysis and knowledge acquisition.

We investigate how to infer the exact relation in a knowledge base for an entity pair from their surface patterns (relations not included in the knowledge base) in the unstructured text.

Idea

We believe that sentences mentioning an entity pair certainly provide some information about the KB relation of this entity pair, though it might not directly describe the relation.

E.g. A triple <Jobs, the-ceo-of, Apple> is in Freebase and a sentence 'Jobs joins Apple' is in Wikipedia. We cannot conclude the former relation from the latter sentence. But if we know more sentences like 'Jobs created Apple' and 'Jobs dominated the development of Apple', we have a higher confidence to infer that Jobs is the ceo of Apple.

Dataset Construction

1. Count the number of pairs of entities for each KB relation in KB (we use Freebase) and sort them. Chose the α (we select 88) commonest relations with corresponding entity pairs. Delete some meaningless entity pairs (e.g. name of the entity is a number.)
2. Extract sentences in Wikipedia containing above entity pairs. If less than β (we select 3) sentences mention an entity pair, this pair is discarded.
3. For each sentence, use Stanford OpenIE to extract the surface pattern between the entity pair.
4. Filter out noise surface patterns by following criterions
 - If a surface pattern(such like 'is in', 'has', 'of', etc.) has more than γ (we select 1000) entity pairs, it is meaningless.
 - Merge surface patterns into a standard pattern with past tense.(e.g. 'go', 'went', 'is going' \rightarrow 'went')
 - Delete surface patterns occurring less than 3 times in Wiki since they have strange semantics.(e.g. <Jobs> *took another medical leave from* <Apple>.)
 - Delete surface patterns labelled with more than 20 KB relations since they are less discriminative.(e.g. <Jobs> *left* <Apple>.)

2. Relation Extraction: From the Perspective of Implicit Relation Reasoning

Model

\mathbf{r} : set of all surface patterns in the training dataset.

r_j : the j th surface pattern in \mathbf{r} .

\mathbf{R} : set of all KB relations in the training dataset.

R_i : the i th KB relation in \mathbf{R} .

$t = \langle e1, e2 \rangle$: an entity pair in the training dataset.

The weight matrix $\boldsymbol{\theta}$ of Dense Layer represents the contribution offered from a surface pattern to a KB relation. (e.g. θ_{ij} quantifies how much contribution provided by r_j to infer the KB relation R_i)

If t has KB relation R_i , then t is an *instance* of R_i , and $\langle t, R_i \rangle$ is an **observed fact**. If t has surface pattern r_j in the training dataset, then t *matches* r_j . If **t is an instance of R_i , and matches r_j** , then **r_j expresses R_i** .

\vec{X} : the **indicator vector of surface patterns** for entity pair t . E.g. $(\vec{X})_j = 1$ if the entity pair t matches r_j , otherwise $(\vec{X})_j = 0$.

To find the contribution matrix $\boldsymbol{\theta}$ that maximizes the likelihood conforming the data, that is:

$$\max_{\boldsymbol{\theta}} \prod_{\langle t, R_i \rangle \in \mathcal{O}_F} P(R_i | t, \mathbf{r}(t), \boldsymbol{\theta})$$

where $P(R_i | t, \mathbf{r}(t), \boldsymbol{\theta})$ is defined as:

$$P(R_i | t, \mathbf{r}(t), \boldsymbol{\theta}) = \frac{\exp(\vec{\theta}_i^T \vec{X})}{\sum_{l=1}^N \exp(\vec{\theta}_l^T \vec{X})}$$

Priori Knowledge Constraints(1)

Observation:

1. Generally, vague surface patterns express more KB relations than surface patterns with specific semantics.

Surface Patterns	KB Relations	Surface Patterns	KB Relations
<i>was in</i>	112	<i>co-founded</i>	2
<i>was with</i>	100	<i>was professor of</i>	2
<i>left</i>	78	<i>was entombed in</i>	1

2. if r_j expresses R_i with more entity pairs, the contribution made by r_j to infer R_i should be **larger**.

Surface Patterns	KB relations	Entity Pairs
<i>spoke</i>	<i>user.original_language</i>	118
<i>university of</i>	<i>education.educational_institution</i>	95
<i>premiered in</i>	<i>film.individual_festivals</i>	78
<i>graduated from</i>	<i>school_type</i>	77
<i>died of</i>	<i>people.cause_of_death</i>	70
<i>college of</i>	<i>education.educational_institution</i>	54

Priori Knowledge Constraints(2)

Constraint 1:

The surface pattern r_j expresses KB relation R_i through some entity pairs, we denotes **the number of these entity pairs** as H_{ij} .

we define $\Theta_{i,j}$ as:

$$\Theta_{i,j} = \frac{H_{i,j}}{\sum_{l=1}^N H_{l,j}}$$

Where the denominator is the total number of entity pairs that match surface pattern r_j in the training dataset. We hope **the greater $\Theta_{i,j}$ is, the greater $\theta_{i,j}$ should be.**

we apply the following constraints on the parameters:

$$(\Theta_{i,j} - \Theta_{i',j'}) (\theta_{i,j} - \theta_{i',j'}) \geq 0, \quad i \neq i' \text{ OR } j \neq j'$$

To **reduce the complexity**, we compare and constrain parameters by fixing surface patterns or KB relations. That is:

$$(\Theta_{i,j} - \Theta_{i',j'}) (\theta_{i,j} - \theta_{i',j'}) \geq 0, \quad i \neq i', j = j' \text{ OR } i = i', j \neq j'$$

Priori Knowledge Constraints(3)

Constraint 2:

Entities have **type** such as {person, location, organization}.

If the types of entity pairs of a KB relation are different from those of a surface pattern, then θ_{ij} should be 0.

There may be multiple types for subject (object) of KB relations and surface patterns. E.g. the types for ***talks about***'s object can be person, organization, location and so on.

Use $S(r_j)$ ($O(r_j)$) to denote the set of appropriate types for the subject (object) of the surface pattern r_j . It is same for the KB relation R_i .

For the surface pattern r_j and KB relation R_i , the type constraints are:

$$\theta_{i,j} = \begin{cases} 0, & S(r_j) \cap S(R_i) = \emptyset \text{ OR } O(r_j) \cap O(R_i) = \emptyset \\ \text{others} \end{cases}$$

Result(1)

Model performance on KBP dataset

We compared our Constrained Model and Simple Model(without prior knowledge) with Baselines: **Mintz ++** (Mintz and Bills 2009), **Hoffmann** (Hoffmann *et al.* 2011), **Surdeanu** (Surdeanu *et al.* 2012), **Han** (Han and Sun 2016).

Model	Precision	Recall	F1 value
Mintz++	0.260	0.250	0.255
Hoffmann	0.306	0.198	0.241
Surdeanu	0.249	0.314	0.278
Han	0.459	0.231	0.307
Simple Model	0.275	0.250	0.262
Constrained Model	0.248	0.449	0.319

Result(2)

Human Evaluation

- We compared our Constrained Model with baselines:(1) **Han**. (2) **Simple model**. (3) **Sebastian** (Riedel *et al.* 2013)
- Set a probability threshold p_0 (we set 0.3) to select the KB relations with conditional probability larger than p_0 . Specifically, for each entity pair t , find 3 KB relations with largest probabilities from the prediction results. Then delete the KB relations with probability less than p_0 .
- We define **the score of t as the sum of the probabilities of selected KB relations**. For each model, we select 100,200,500 entity pairs with the highest score as evaluation objects.
- We assigned the truth or false of each relation according to **the majority vote of the labels by three Ph.D. students**.

Precision (%)	Top 100	Top 200	Top 500	Average
Han	0.86	0.83	0.74	0.81
Sebastian	0.80	0.76	0.63	0.73
Simple model	0.77	0.73	0.62	0.71
Constrained model	0.89	0.85	0.79	0.84

2. Relation Extraction: From the Perspective of Implicit Relation Reasoning

Result(3)

Typical Case

It is easy for our human to find that those Surface Patterns contribute to inferring the KB Relations at the left. It proves that our method works well.

KB Relations	Surface Patterns (Contribution Weights)
<i>tv_program.languages</i>	<i>was described in</i> (1.33), <i>with the language of</i> (3.21), <i>was screened in</i> (2.13)
<i>people.place_of_birth</i>	<i>wrote poetry in</i> (0.31), <i>was working in</i> (1.87), <i>taught English in</i> (0.25), <i>was born in</i> (2.83), <i>grow up in</i> (1.23)
<i>people.cause_of_death</i>	<i>died of</i> (3.12), <i>cause of death was</i> (2.13), <i>died from</i> (1.09), <i>was diagnosed with</i> (0.65), <i>shortly succumbed to</i> (0.89)
<i>person.nationality</i>	<i>returned once again in</i> (0.91), <i>was working in</i> (1.32), <i>voted for</i> (1.78), <i>eventually died in</i> (0.79), <i>next king of</i> (0.68)
<i>people.profession</i>	<i>was fond of</i> (1.36), <i>rely on</i> (0.87), <i>was working as</i> (4.31), <i>was one of the</i> (3.11), <i>was a</i> (1.20), <i>acted on</i> (1.88)
<i>film.directed_by</i>	<i>was produced by</i> (1.39), <i>the baby of</i> (0.52), <i>was elaborated by</i> (5.89) <i>directed by</i> (7.31), <i>the director's name</i> (4.27)
<i>people.sibling_relationship</i>	<i>loved</i> (1.23), <i>the darling of</i> (0.98), <i>successor of</i> (1.25), <i>accompanied</i> (3.11) <i>lived with</i> (2.45)
<i>award.honored_for</i>	<i>was famous for</i> (1.02), <i>had a talent in</i> (2.19), <i>was enjoying</i> (1.09), <i>was working as</i> (4.21), <i>was interested in</i> (2.39), <i>was hero in</i> (0.79)
<i>tv.tv_actor</i>	<i>was a star in</i> (1.36), <i>liked</i> (1.09), <i>was busy making</i> (3.10), <i>was famous in</i> (3.18), <i>a actor in</i> (4.14), <i>a famous actor in</i> (5.21)

- 3. Hierarchical Conceptual Labelling

Problem

we usually adopt the bag of words (BoW) model to represent the meaning of the text. It is usually a collection of scattered words, so it is difficult for human beings or machines to understand when its size is large. It is very important to explicitly represent the semantics of a BoW.

Idea

we focus on **conceptualizing** a BoW. The use of concepts as labels, known as **conceptual labeling**, makes a good semantic representation of a set of words. We use **Bayesian Rose Trees** to hierarchically cluster.

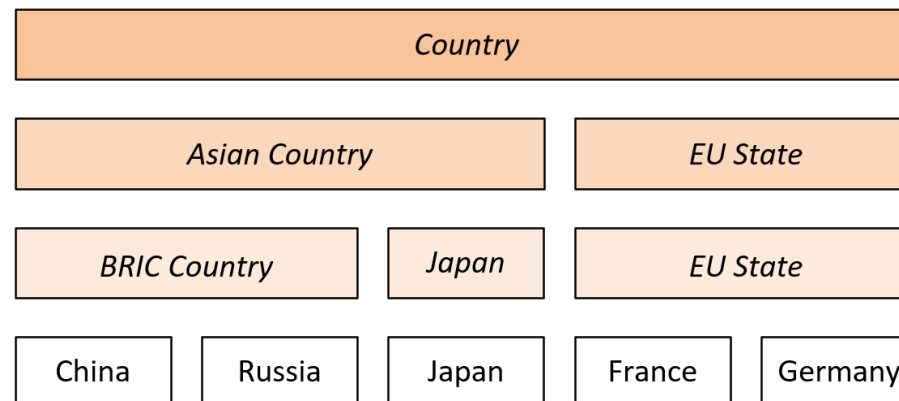
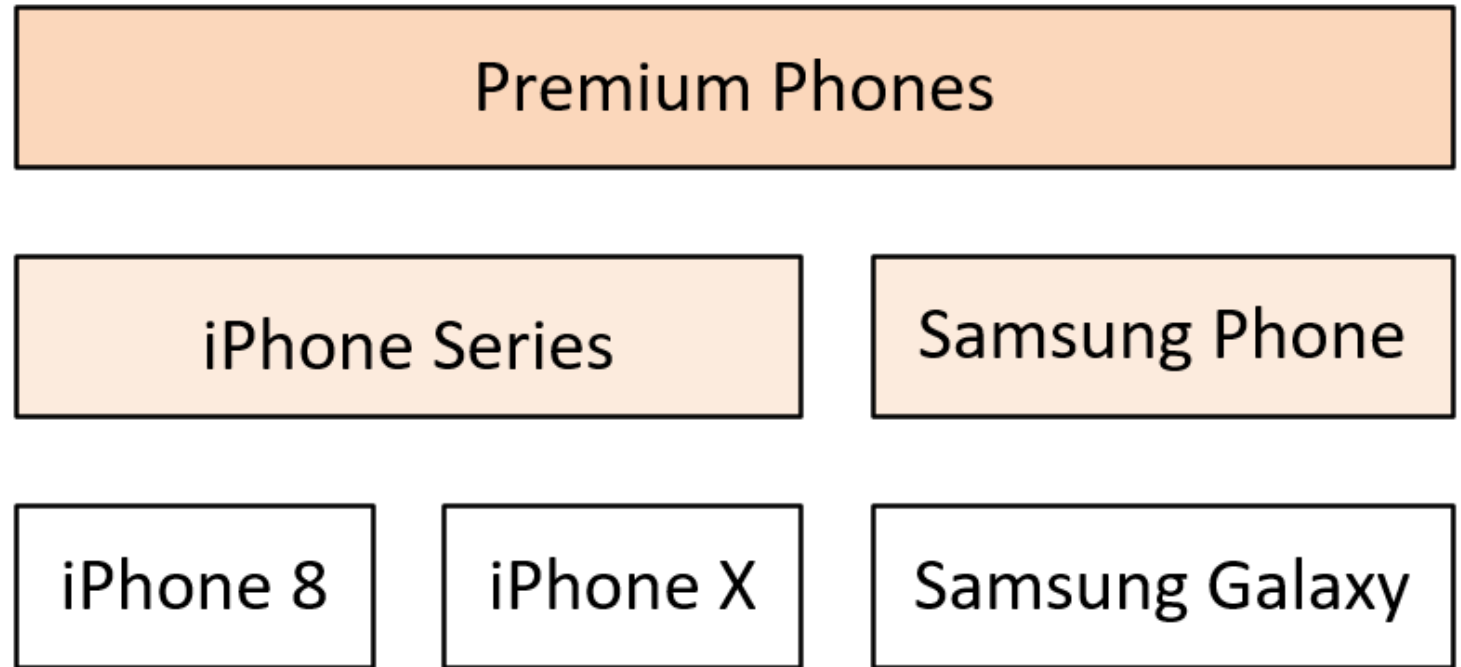


Figure: Conceptualize a Bow Hierarchically

Application

Method proposed in our paper have various applications. Especially in **interactive search**, we can use historical data to make **accurate recommendation**. E.g. If we search 'iPhone 8' at first, the results are iPhone 8 related. Then if we search 'iPhone X', our hierarchical concept could be 'iPhone series phone'. If we input 'Samsung Galaxy' next time, our concept could become 'premium phones' **based on the Bow 'iPhone 8, iPhone X, Samsung Galaxy'**.

Using above concept, we could improve our recommendation system greatly.



Using Probase Knowledge

Probase provides lots of **concept entity** pairs in the form of **<concept, entity, frequency>** we introduce the **typicality** score which is defined as (Song *et al.* 2015):

$$p(e|c) = \frac{n(c, e)}{\sum_{e_i} n(c, e_i)} \quad p(c|e) = \frac{n(c, e)}{\sum_{c_i} n(c_i, e)}$$

Define a prior probability of a concept or an entity :

$$p(c) = \frac{\sum_{e_i} n(c, e_i)}{\sum_{(c_j, e_i)} n(c_j, e_i)} \quad p(e) = \frac{\sum_{c_j} n(c_j, e)}{\sum_{(c_j, e_i)} n(c_j, e_i)}$$

Filter out noise

e.g. given the BoW {university, China, Japan, U.S.A}, it is clear that all the words are countries except **university**, so **University** is likely to be noise.

Basic feature of noisy word: **hard to be merged with any other word**.

let D be the input BoW. d_i and d_j be the i th and j th entity (word) in D , respectively.

We use $p(c|d_i, d_j)$ to measure **the degree that the concept c summarizes the semantics of the two entities d_i, d_j** .

$$\begin{aligned} p(c|d_i, d_j) &= \frac{p(d_i, d_j|c)p(c)}{p(d_i, d_j)} \\ &= \frac{p(d_i|c)p(d_j|c)p(c)}{p(d_i)p(d_j)} \end{aligned}$$

Assuming that all entities in D have same probability, $p(d_i) = \tilde{p}$.

Introduce a hyper-parameter δ to control the "appropriate" degree that we regard d_i as noise if it satisfies that:

$$\max_{c \in \mathcal{C}_{i,j}, d_j \in \mathcal{D}} p(c|d_i, d_j) < \delta$$

Bayesian Rose Trees

BRT(Bayesian Rose Trees) is a novel hierarchical clustering algorithm that **builds a tree from bottom up** by combining the **two most similar** clusters at each step.

To determine **which T_i and T_j should be merged as well as which operation should be selected**, BRT takes the criterion of maximizing the likelihood ratio $L(T_m)$, which is defined as follows:

$$L(T_m) = \frac{p(\mathcal{D}_m | T_m)}{p(\mathcal{D}_i | T_i) p(\mathcal{D}_j | T_j)}$$

Where $p(\mathcal{D}_m | T_m)$ is the likelihood of data \mathcal{D}_m given the tree T_m , which is recursively calculated as follows:

$$p(\mathcal{D}_m | T_m) = \pi_m f(\mathcal{D}_m) + (1 - \pi_m) \prod_{T_k \in \text{ch}(T_m)} p(\mathcal{D}_k | T_k)$$

$f(\mathcal{D}_m)$ is the marginal probability of data \mathcal{D}_m . Intuitively, π_m is the prior probability that all the data in T_m is kept in one cluster instead of partitioned into sub-trees. $\text{ch}(T_m)$ is the abbreviation of $\text{children}(T_m)$.

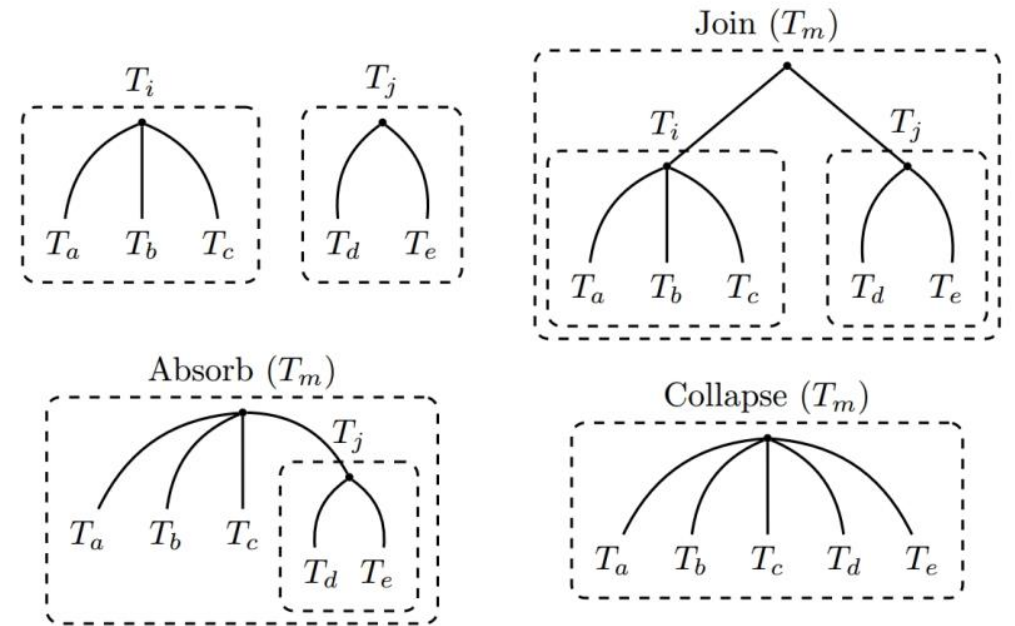


Figure: Three possible merging operations in BRT

Hierarchical Conceptual Labelling(1)

- $f(D_m)$ denotes the probability that all the data points in D_m are generated by the same probability model. D_m can be considered to be generated by concepts in Probase. Define C_m be the common concepts set of D_m . E.g. let D_m be {China,Brazil,India}, its C_m is {*developing country, emerging market, BRIC country*}.

We define the **selected probability for c_i** , that is:

$$p_s(c_i) = \frac{p(c_i)}{\sum_{c \in C_m} p(c)}$$

Then $f(D_m)$ is computed as follows:

$$f(D_m) = \sum_{c \in C_m} p_s(c) p(D_m | c)$$

- We set the prior $\pi_m = 0.5$ since **we have no prior knowledge** to determine which partitions are more important.

Hierarchical Conceptual Labelling(2)

- The following criterion is used to select **the most appropriate conceptual label**:

$$\begin{aligned} c_m^* &= \arg \max_{c \in \mathcal{C}_m} p(c | \mathcal{D}_m) \\ &= \arg \max_{c \in \mathcal{C}_m} p(\mathcal{D}_m | c) p(c) \end{aligned}$$

- We should stop clustering **when there is no appropriate label that well summarize the current cluster**. E.g. if $\mathcal{D}_m = \{\text{children, cat, cow, elephant}\}$, we can only use labels, such as *creature*, to summarize it, which makes no sense since they are **too vague**. To avoid it, we introduce the likelihood ratio threshold γ and stop clustering when $L(T_m) < \gamma$.

Algorithm

- D is the clean BoW whose noise has been filtered out.
- In the initialization, we set $p(D_i/T_i) = 1$ ($i = 1, \dots, N$) and $L(T_m) = \gamma_0$ ($\gamma_0 > \gamma$).
- In each step of clustering, the selection of the pair of trees T_i and T_j as well as the merge operation are determined by Equation 13.
- We select the concept in C_m that maximizes $p(c/D_m)$ as the label of the cluster D_m (see Equation 14).
- The procedure is repeated until $L(T_m) < \gamma$

Algorithm 1 Hierarchical conceptual labeling based on Bayesian rose tree

Input:

data $\mathcal{D} = \{d_1, d_2, \dots, d_N\}$ and $LabelSet = \{\}$

Initialize:

number of clusters $k = N$;

$\mathcal{D}_i = \{d_i\}$, $T_i = \{d_i\}$ and $p(\mathcal{D}_i/T_i) = 1$ ($i = 1, \dots, N$);

$L(T_m) = \gamma_0$

while $k > 1$ **and** $L(T_m) > \gamma$ **do**

Find the pair of trees T_i and T_j and the merge operation that maximize the likelihood ratio:

$$L(T_m) = \frac{p(\mathcal{D}_m | T_m)}{p(\mathcal{D}_i | T_i) p(\mathcal{D}_j | T_j)} \quad (13)$$

Select the conceptual label c_m^* :

$$c_m^* = \arg \max_{c \in C_m} p(c | \mathcal{D}_m) \quad (14)$$

Merge T_i and T_j into T_m using the selected merge operation;

$\mathcal{D}_m \leftarrow \mathcal{D}_i \cup \mathcal{D}_j$;

Add c_m to $LabelSet$;

Delete T_i and T_j , $k \leftarrow k - 1$

end while

Result on Synthetic Data—Generation

Two ways to generate Synthetic data:

1. For each BoW, selecting m concepts from Probase, then randomly select n instances from these m concepts' entities (ensure that at least two entities were selected in each concept), so we have a BoW with size n . Noise words are added to the BoW which are randomly selected from other concepts' entities.
2. For each Bow, select m concepts randomly in Probase and then n entities are selected randomly in each concept and l noise words are selected, then we get a BoW with size $mn+l$.

Result on Synthetic Data—Denoising

- We choose the first data generation method.
- Constitute 5 BoWs with $m=2$ and $n=5$, $m=3$ and $n=10$, $m=5$ and $n=15$ (20, 25). Add 10% noise words (**noted as Set_B , are randomly selected from other concepts' entities**). Use denoising algorithm ($\delta_{noise} = 5 \times 10^{-8}$) to get a **new noise word set Set_A** .
- We denote the union of difference Set_A and difference Set_B as M_{noise} , the size of the BoW (including noise words) is N_{BoW} , then the **Accuracy** of the denoising algorithm is defined as $(1 - M_{noise})/N_{BoW}$

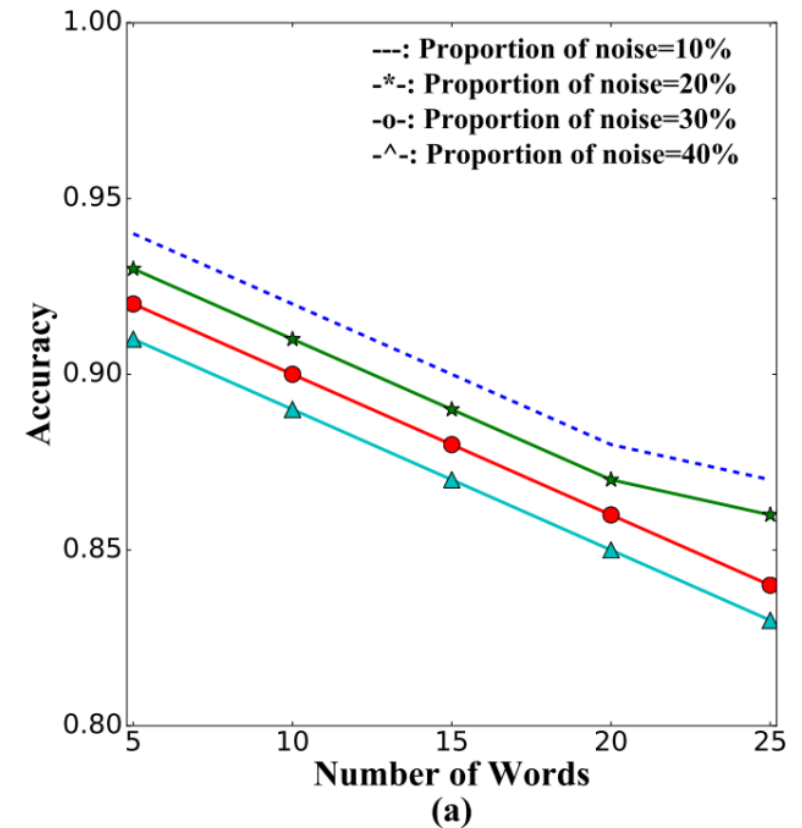


Figure: Results of denoising

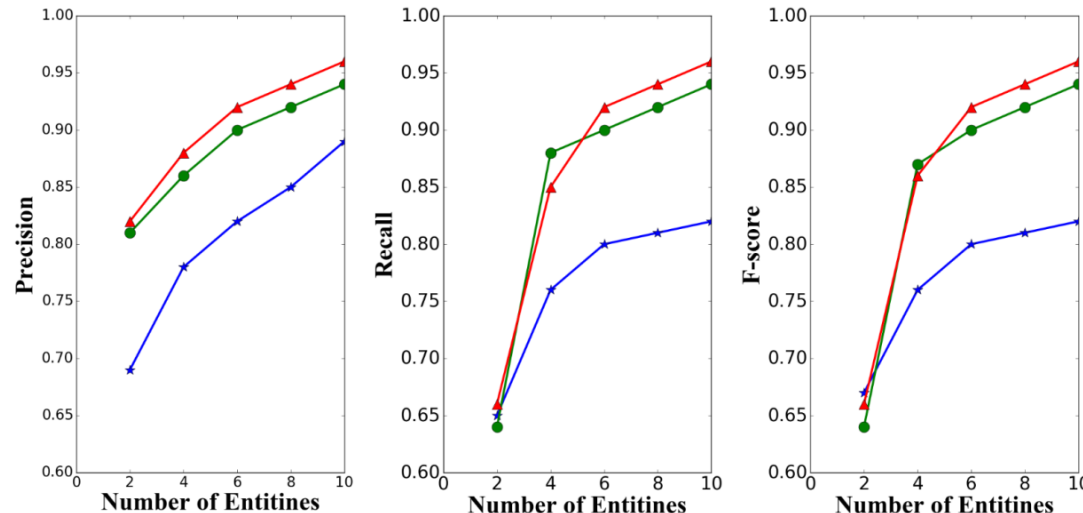
Result on Synthetic Data—Labelling

- Compare our method with the method for generating conceptual labels based on **MDL** principle and **CC**.
- Choose the **second** data generation method.
- Set $L(T_m)$ threshold γ as 0.8. Selected **conceptual labels in the top layer** and compared them with **pre-given 5 concepts**.
- Assuming for the BoW i , the size of conceptual labels in the top layer is x_i , and y_i of them are in the pre-given label set(5 concepts).

Define: Precision: $\left(\sum_{i=1}^t y_i \right) / \left(\sum_{i=1}^t x_i \right)$

Recall: $\left(\sum_{i=1}^t y_i \right) / 5t$

F: $\frac{2precision \times recall}{precision + recall}$



Curve with triangle: BRT
Curve with circle: MDL
Curve with star: CC

Result on Synthetic Data—Human Evaluation

- N is the number of labels.
- We choose the **first** data generation method.
- Evaluation Standard:

Score	Description	Example
$1/N$	Appropriate	football,sport,snooker→sport
$1/2N$	General	art,painting,dance,disco→activity
0	Inappropriate	family,sibling,people→others

- Evaluate the performance (accuracy) of the algorithm by calculating the scores of the all conceptual labels.

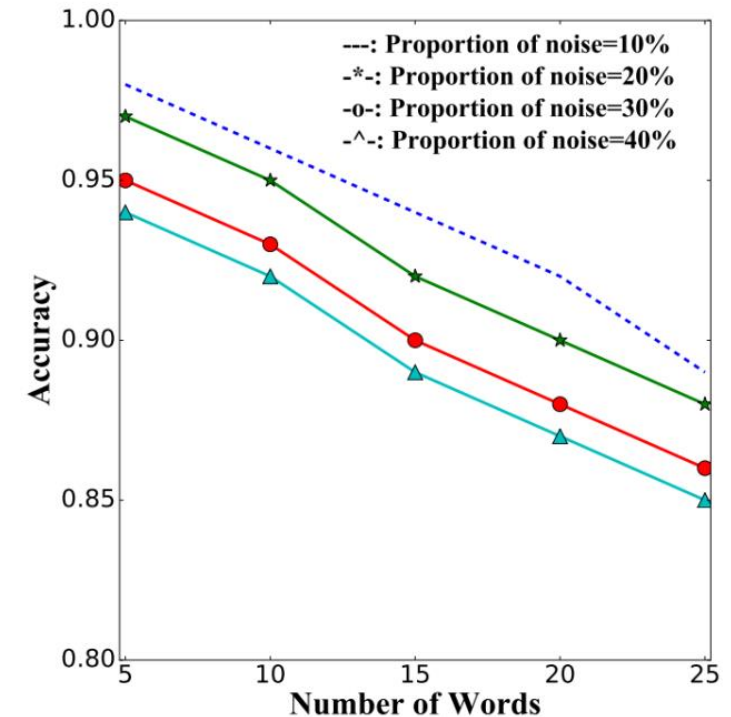


Figure: Results with different noise ratios and sizes of BOW

Result on Real Data—Human Evaluation

- Use Wiki data and Flickr data. Use LDA algorithm and extracted top words of each topic as the BoW.
- Use manual scoring method to evaluate the performance of the conceptual label.
- Evaluation scores on Flickr and Wikipedia data:

Algorithm	Flickr data	Wikipedia data
CC	0.76	0.81
MDL Method	0.85	0.88
BRT Method	0.91	0.94

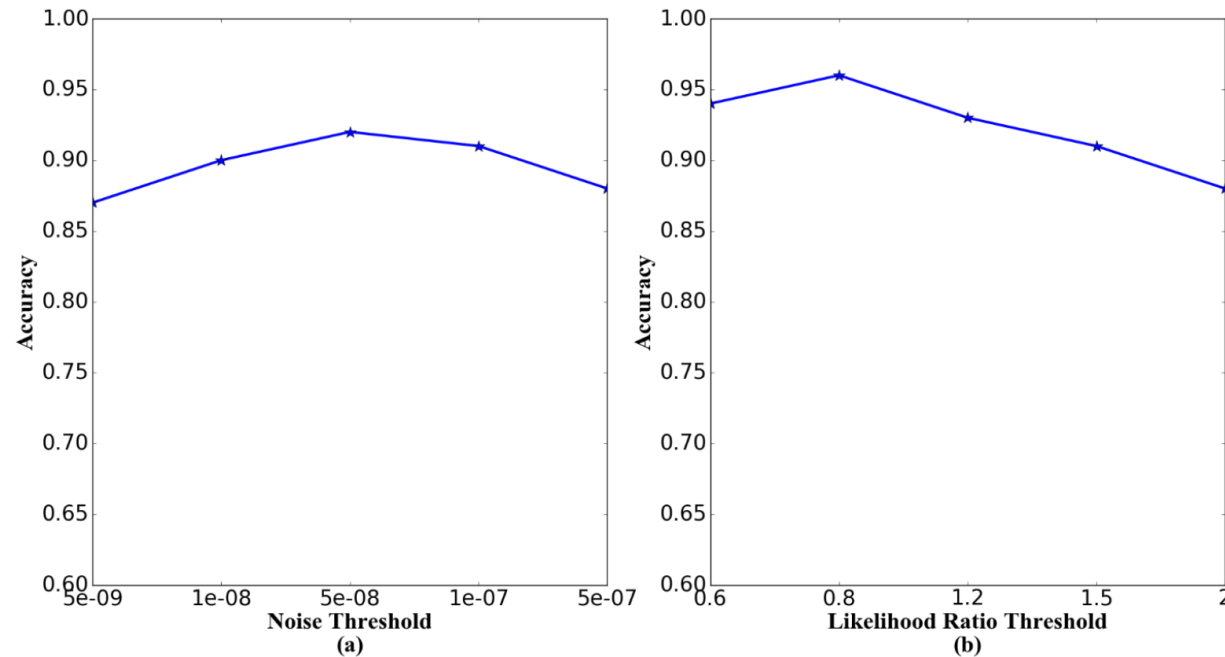
Two empirical threshold

There are two empirical threshold in our Algorithm:

1. **Noise threshold** δ_{noise} .
2. **Likelihood ratio threshold** $L(T_m) \gamma$

We set the size of the BoW is 10, the noise ratio is 20%, the effect of different δ_{noise} is analyzed in the figure(a):

Based on the above analysis, we set $\delta_{noise} = 5 \times 10^{-8}$ and analyse the likelihood ratio threshold $L(T_m) \gamma$ in figure(b):

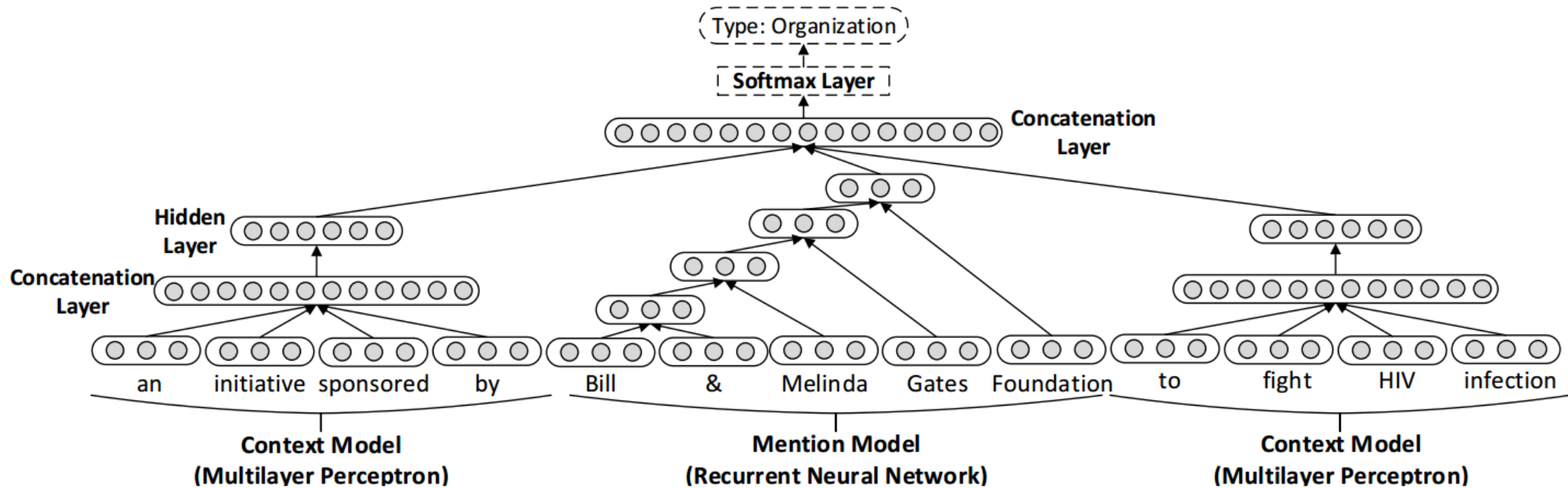


- 4. Entity Typing with DNN and Word Embedding

*Follow Li Dong, Fu-Ru Wei, etc. A Hybrid Neural Model for Type Classification of Entity Mentions. In IJCAI, 2015.

Baseline Model

*Li Dong, Fu-Ru Wei, etc. A Hybrid Neural Model for Type Classification of Entity Mentions. In IJCAI, 2015.



It is impressive for this model to consider Context Model. However we believe that using Multilayer Perceptron and Word Embedding is not enough.

Character-level Word Embedding

*Improved from Yanyao Shen, Hyokun Yun, etc. Deep Active Learning for Named Entity Recognition. On arxiv.

- We use Bidirectional LSTM to learn more knowledge than CNN in Shen's paper.
- Character-level Word Embedding can be a good supplement to typical Word Embedding.

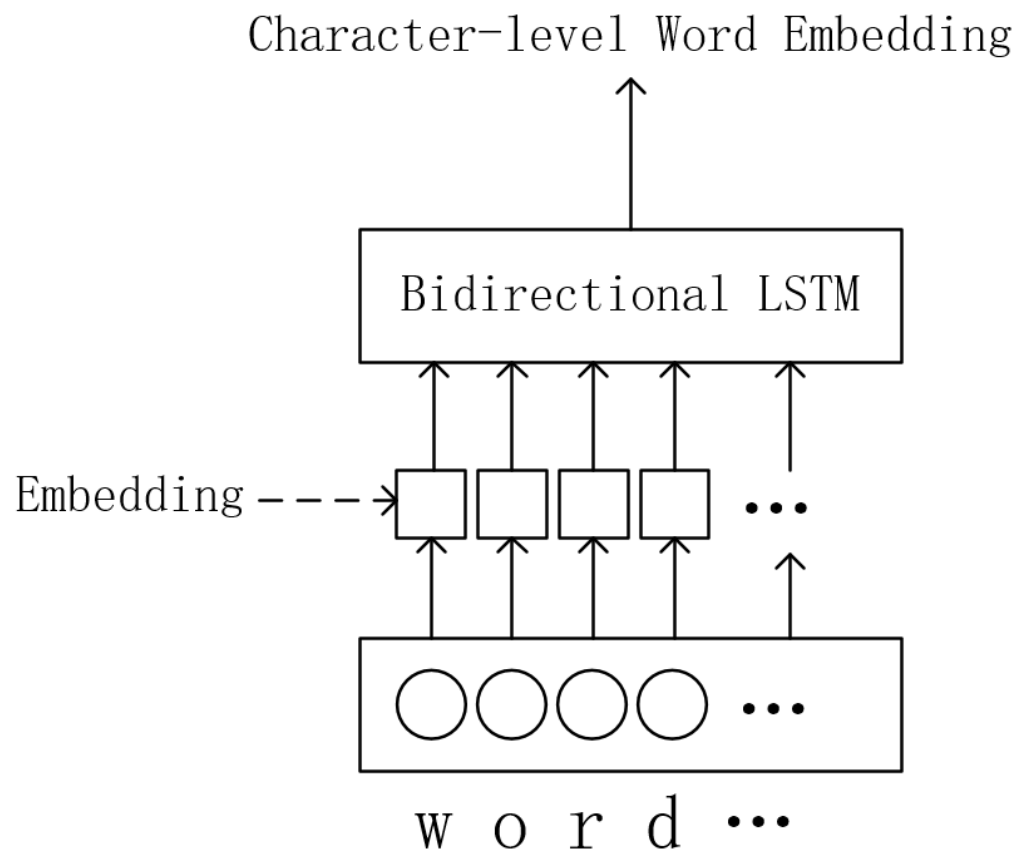
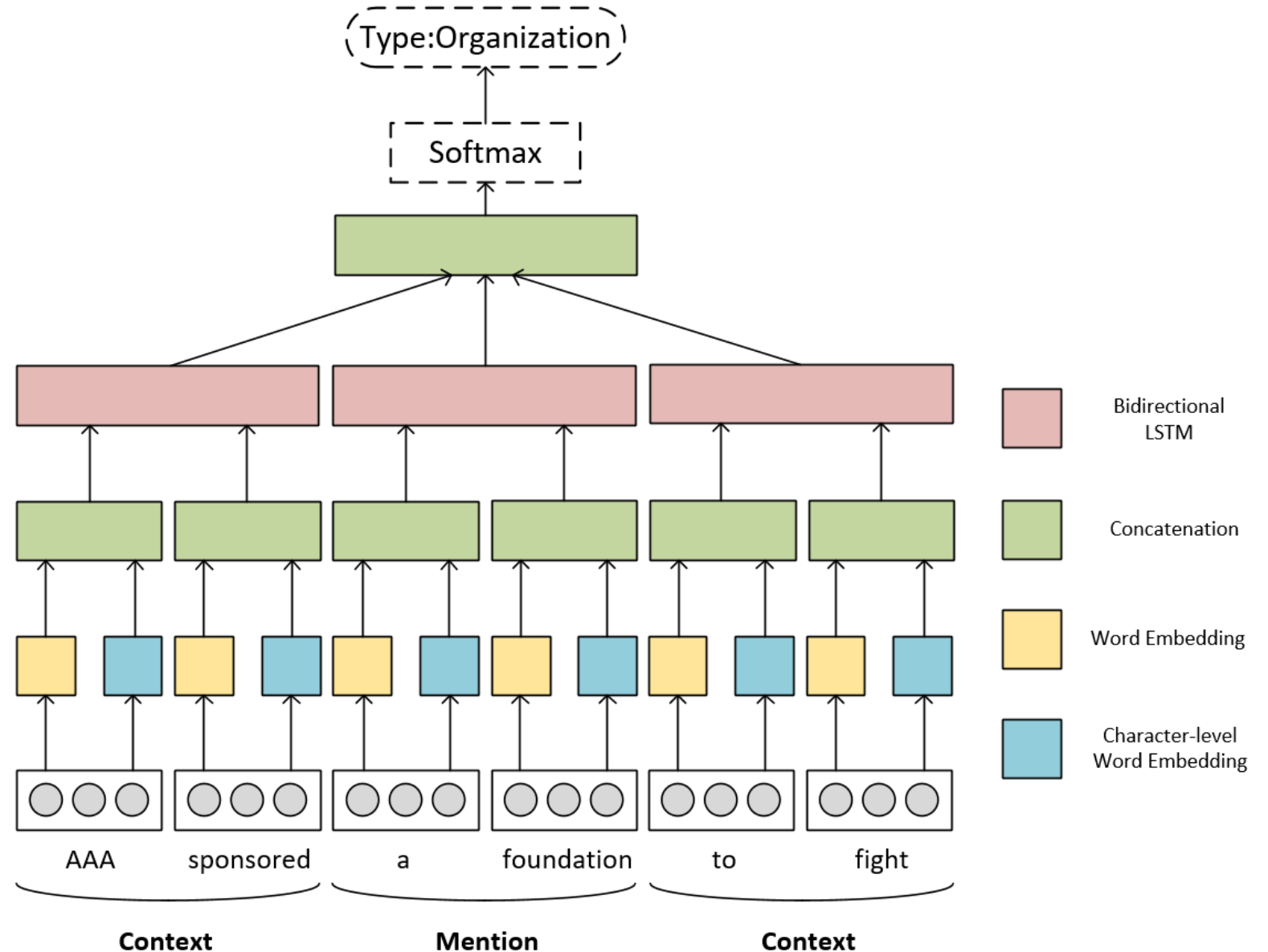


Figure: Character-level Word Embedding structure

Model

- Replace Multilayer Conception with Bidirectional LSTM.
- Add improved Character-level Word Embedding.
- The length of Context and Mention Windows can be adjusted. We choose Mention Window = 6 and Context Window = 7 to balance accuracy and training speed.



4. Entity Typing with DNN and Word Embedding

Result

- Do experiment on dataset Wiki, BBN, and OntoNotes.
- The results are as follows.

	Wiki		OntoNotes		BBN	
	Ma-F1	Mi-F1	Ma-F1	Mi-F1	Ma-F1	Mi-F1
Baseline	0.409	0.417	0.288	0.272	0.591	0.606
Our model	0.602	0.663	0.454	0.421	0.632	0.658

Thanks !

If you have any questions or interests, please contact me without hesitation. ☺

Zhe Xu

Shanghai Key Laboratory of Data Science

Fudan University

Email: zxu14@fudan.edu.cn