

Deep Learning for Natural Language Understanding

GU JIATAO

DEPARTMENT OF ELECTRICAL AND ELECTRONIC ENGINEERING

THE UNIVERSITY OF HONG KONG

JIATAOGU@EEE.HKU.HK

Natural Language Understanding (NLU)

Understanding human language is difficult to represent:

- Syntax
- Semantics
- Pragmatics

Applications:

- Topic Modelling
- Information Extraction
- Machine Translation
- Text Summarization
- Dialogue Systems
- ...

Representation Learning for NLU

Representation Learning:

- is concerned with questions surrounding how we can best **learn** meaningful and useful representations of data.
- The performance of machine learning methods is heavily dependent on the choice of data representation on which they are applied.

Deep Learning:

- Use deep learning (deep neural network) methods to generate proper data representation
- Typically, **vectorial/distributed representation**.

Undirected Topic Models

Modelling a document using a Restrict Boltzmann Machine (RBM).

- **Traditional Learning Algorithm:**
 - Contrastive Divergence (CD):
 - slow and unstable
- **Our methods:**
 - Noise Contrastive Estimation
 - Fast but cannot directly applied to documents.
 - α -NCE:
 - **modifying the traditional NCE and applicable to natural language**
 - **much faster than CD.**
 - **performance comparable**

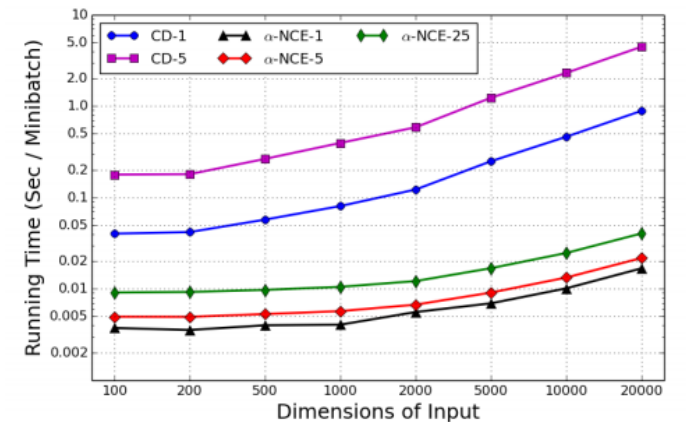
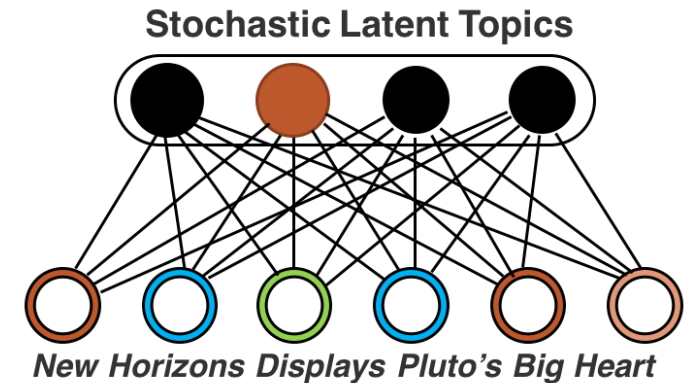


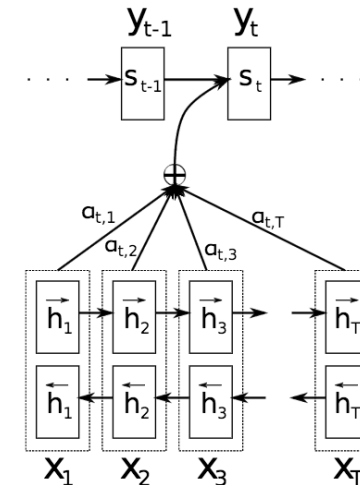
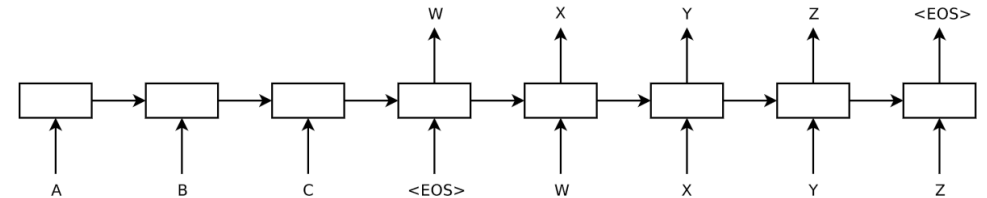
Figure 1: Comparison of running time

Sequence-to-Sequence Learning

Different from Topic Modelling, this part focus more on word orders and has more applications on real-life, such as machine translation.

The core idea is “**Sequence-to-sequence Learning**”:

- We have:
 - Source sequence
 - Target sequence
- What we do:
 - Model **the transformation** using a **Encoder** and **Decoder**
 - With additional techniques such as **Attention Mechanism**
- **What we can do more:**
 - **Look deep into the Seq2Seq learning framework.**



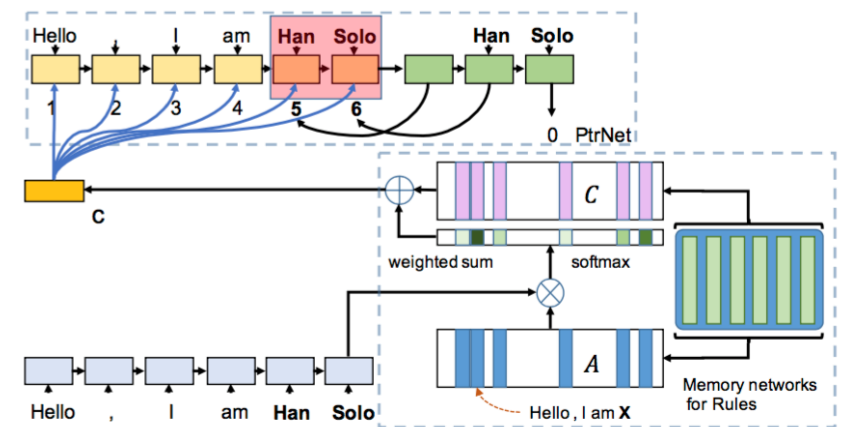
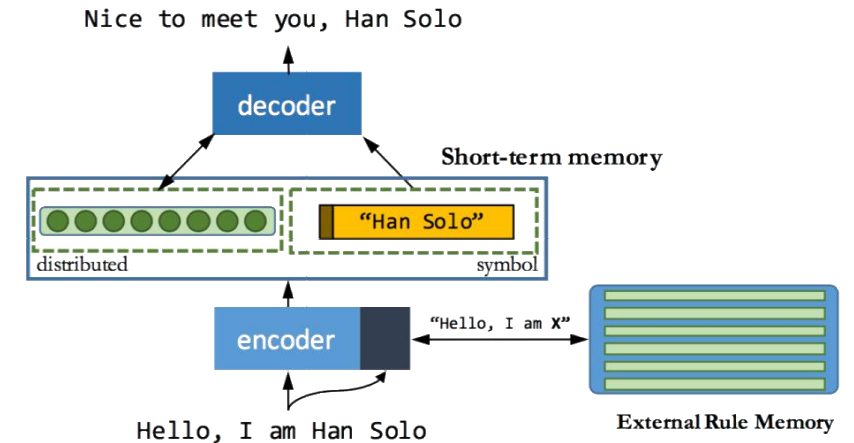
Incorporating Transformation Rules

Basic Idea:

- Incorporating rules into Seq2Seq Learning

Useful in dialogue system:

- A simple dialogue system can be seen as a **Seq2Seq Learning framework**
- Transformation Rules:
 - *Hello, I am X -> Nice to meet you, X*
- Combine existing transformation rules into a neural network based dialogue system.



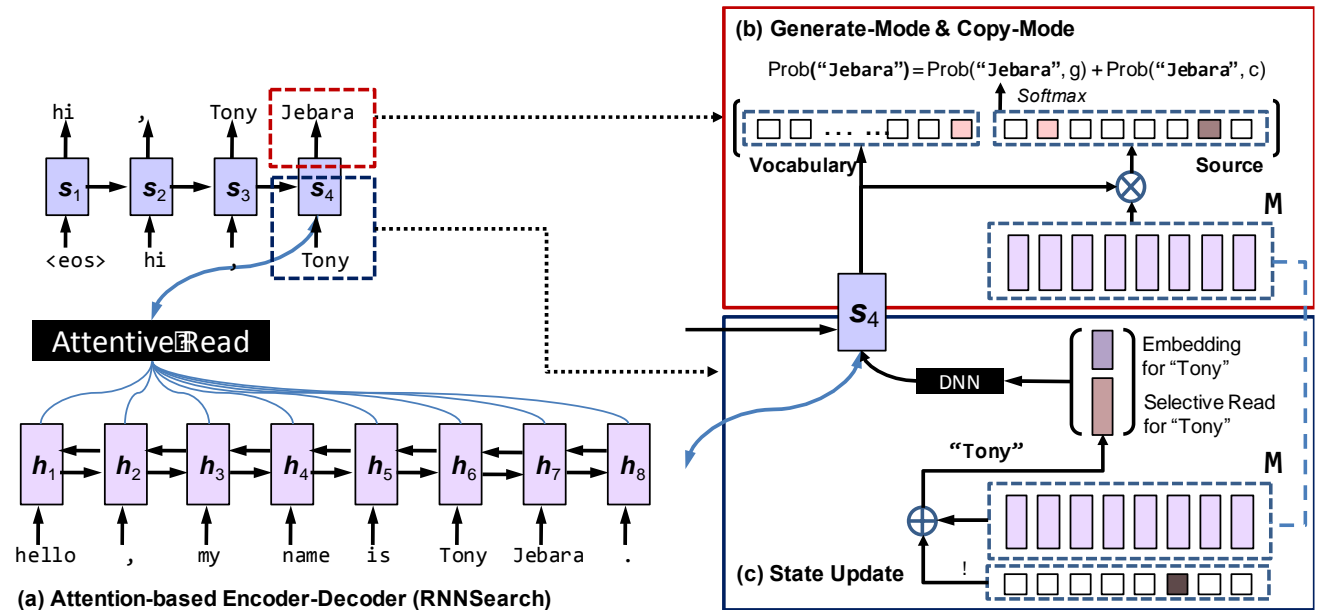
Incorporating Copying Mechanism

Basic Idea:

- Incorporating a special Copying Mechanism in traditional Seq2Seq learning framework.

Useful in **text summarization, machine translation and dialogue system**:

- Sometimes we not only need generate words, but also need copy from the source.
- **A hybrid addressing mechanism of both context and location.**



Incorporating Copying Mechanism

Input(3): 今天上午9点半, 复旦投毒案将在上海二中院公开审理。被害学生黄洋的亲属已从四川抵达上海, 其父称待刑事部分结束后, 再提民事赔偿, 黄洋92岁的奶奶依然不知情。今年4月, 在复旦上海医学院读研究生的黄洋疑遭室友林森浩投毒, 不幸身亡。新民网
Today 9:30, the Fudan poisoning case will be will on public trial at the Shanghai Second Intermediate Court. The relatives of the murdered student Huang Yang has arrived at Shanghai from Sichuan. His father said that they will start the lawsuit for civil compensation after the criminal section. HuangYang 92-year-old grandmother is still unaware of his death. In April, a graduate student at Fudan University Shanghai Medical College, Huang Yang is allegedly poisoned and killed by his roommate Lin Senhao. Reported by Xinmin

Golden: 林森浩投毒案今日开审92岁奶奶尚不知情 (the case of Lin Senhao poisoning is on trial today, his 92-year-old grandmother is still unaware of this)

RNN context: 复旦投毒案: 黄洋疑遭室友投毒凶手已从四川飞往上海, 父亲命案另有4人被通知家属不治?

CopyNet: 复旦投毒案今在沪上公开审理 (the Fudan poisoning case is on public trial today in Shanghai)

Input(4): 华谊兄弟 (300027) 在昨日收盘后发布公告称, 公司拟以自有资金3.978亿元收购浙江永乐影视股份有限公司若干股东持有的永乐影视51%的股权。对于此项收购, 华谊兄弟董秘胡明昨日表示: “和永乐影视的合并是对华谊兄弟电视剧业务的一个加强。
Huayi Brothers (300027) announced that the company intends to buy with its own fund 397.8 million 51% of Zhejiang Yongle Film LTD's stake owned by a number of shareholders of Yongle Film LTD. For this acquisition, the secretary of the board, Hu Ming, said yesterday: "the merging with Yongle Film is to strengthen Huayi Brothers on TV business".

Golden: 华谊兄弟拟收购永乐影视51%股权 (Huayi Brothers intends to acquire 51% stake of Zhejiang Yongle Film)

RNN context: 华谊兄弟收购永乐影视51%股权: 与永乐影视合并为“和唐”影视合并的“UNK”和“UNK”的区别?

CopyNet: 华谊兄弟拟3.978亿收购永乐影视董秘称加强电视剧业务 (Huayi Brothers is intended to 397.8 million acquisition of Yongle Film secretaries called to strengthen the TV business)

Input(7): 工厂, 大门紧锁, 约20名工人散坐在树荫下。“我们就是普通工人, 在这里等工资。”其中一人说道。7月4日上午, 记者抵达深圳龙华区清湖路上的深圳愿景光电子有限公司。正如传言一般, 愿景光电子倒闭了, 大股东邢毅不知所踪。
The factory's door is locked. About 20 works are scattered to sit under the shade. "We are ordinary workers, we are waiting for our salary here." one of them said. In the morning of July 4th, reporters arrived at Shenzhen Yuanjing Photoelectron Corporation located at Qinghu Road, Longhua District, Shenzhen. Just as the rumor says, Yuanjing Photoelectron Corporation is closed down and the large shareholder Xing Yi is missing.

Golden: 深圳亿元级LED企业倒闭烈日下工人苦等老板 (Hundred-million CNY worth LED enterprise is closed down and workers wait for the boss under the scorching sun)

RNN context: 深圳“<UNK>”: 深圳<UNK><UNK>, <UNK>, <UNK>, <UNK>

CopyNet: 愿景光电子倒闭20名工人散坐在树荫下 (Yuanjing Photoelectron Corporation is closed down, 20 works are scattered to sit under the shade)