

FOUNDATION

- Supervised Learning (Lecture 1-5)
 - Not Efficient to collect data for each task
- Self-supervised Learning (Lecture 7)
 - Pre-train (Develop general purpose knowledge)
 - BERT GPT-3 GPT-2
 - ELMO T5
- Generative Adversarial Network (Lecture 6)
- Reinforcement Learning (Lecture 12)

Model Training

Function with Unknown Parameters

Domain Knowledge 知识域知识

Define Loss from Training Data

- Loss 关于参数的函数 $L(b, w)$
衡量一组参数的好坏

- MAE (Mean Absolute Error) $e = |y - \hat{y}|$

- MSE (Mean Square Error) $e = (y - \hat{y})^2$

y 和 \hat{y} 都是概率 \rightarrow Cross-entropy

Optimization $w^*, b^* = \arg \min L$

- Hyperparameters 自己定义的超参数 学习率

- Local minima 不是 Gradient descent 真正的问题

REGRESSION

Step 1: function with unknown

Linear Function

Model Bias 由模型导致的误差

More Flexible \rightarrow 单个特征

多个特征

Step 1:
function with unknown

Step 2: define loss

Step 3:
optimization

- 参数越多，越容易 overfit
- 未知参数的可能性(由选择的 Function 放置) \rightarrow Complexity

We want $L(h^{\text{train}}, D_{\text{all}}) - L(h^{\text{all}}, D_{\text{all}}) \leq \delta$

What kind of D_{train} fulfill it?

We Need

$$\forall h \in \mathcal{H}, |L(h, D_{\text{train}}) - L(h, D_{\text{all}})| \leq \delta/2$$

D_{train} is a good proxy of D_{all} for evaluating loss L given any h . 使理想接近现实的 D_{train}

Formulation:
 $L(h^{\text{train}}, D_{\text{all}}) \leq L(h^{\text{train}}, D_{\text{train}}) + \delta/2$

$$\leq L(h^{\text{all}}, D_{\text{train}}) + \delta/2$$

$$\leq L(h^{\text{all}}, D_{\text{all}}) + \delta/2 + \delta/2 = L(h^{\text{all}}, D_{\text{all}}) + \delta$$

$P(D_{\text{train}} \text{ is bad}) \leq |\mathcal{H}| \cdot 2 \exp(-2N\varepsilon^2)$ N 越大， M 越小 \rightarrow Sample 到差 D_{train} 的概率

N : D_{train} 中的样本数量 \rightarrow 训练集规模 \uparrow , 模型复杂度 \downarrow If we want $P(D_{\text{train}} \text{ is bad}) \leq \delta$

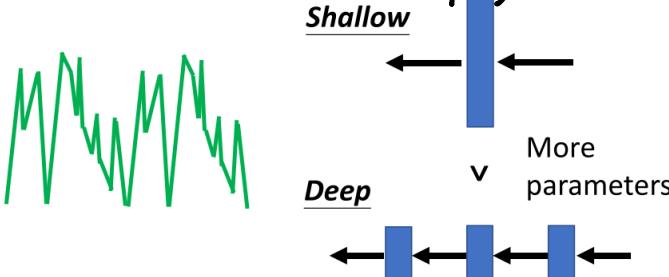
$|\mathcal{H}|$: 由选择的 Function 放置 (Model Complexity)

(ε : 自定义参数 1 越小, 对 D_{train} 的标准越高) \rightarrow 确定训练集规模

When parameters are continuous: VC-Dimension 判断 Complexity

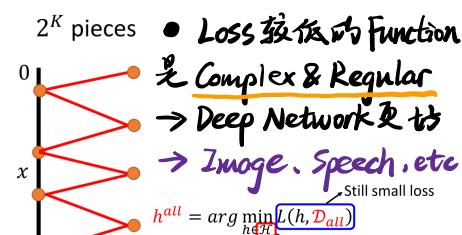
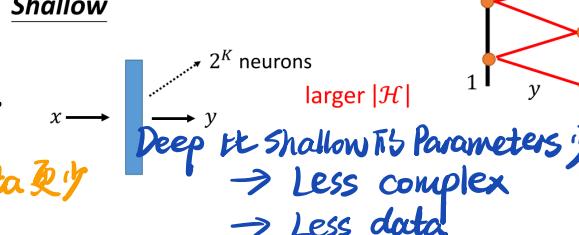
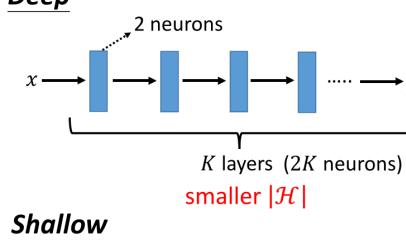
Dilemma: 很难采集到较大的 N

样本过少, h 选择太少, 效果差 $\rightarrow L(h^{\text{train}}, D_{\text{all}})$ 与 $L(h^{\text{all}}, D_{\text{all}})$ 相差很小, 但 $L(h^{\text{all}}, D_{\text{all}})$ 本身很大



• 一层 Hidden Layer 可以表示任何 Function
但 Deep Structure 更有效

\rightarrow 不易 Overfitting, 需要的数据更少



• Deep \rightarrow 更低的 Loss
 \rightarrow 更少的 Candidates

CNN

假设所有输入的 image 大小相同

\rightarrow 先 Rescale 成一样大小

Receptive Field 感受野 \rightarrow 每个 Neuron 只看局部的 Pattern

- 不同 Neuron 可以重叠, 成 Receptive Field 相同
- 可以只有部分 channel 每个 Receptive Field 一般会有多个 Neuron (64 or 128)
- 形状可以自定

Kernel 卷积核

Receptive Field 的长 x 宽

Stride 步长 (Hyper Parameter)
 \rightarrow RF 移动的范围 (1 or 2, 一般有 overlap)

Padding 填充

\rightarrow 超出范围后向 RF 中填充

Parameter Sharing 参数共享

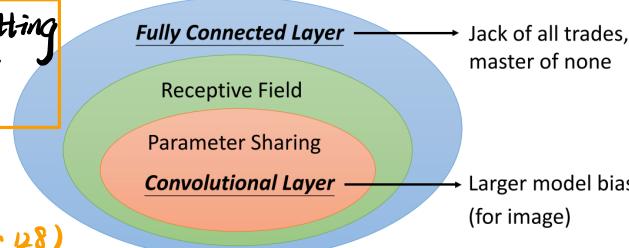
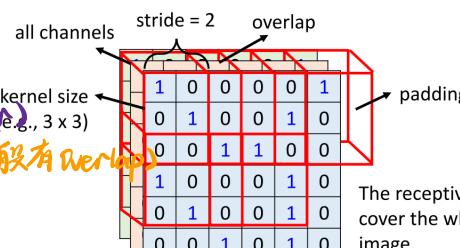
- 相同的 Pattern 可能出现在图片的不同区域 (break)

\rightarrow Parameter Sharing

- 两组 RF 不同的 Neuron 的 weight 相同 (权值共享) \rightarrow RF 相同的 Neuron 不会共享参数

\rightarrow 共用的参数 \rightarrow Filter

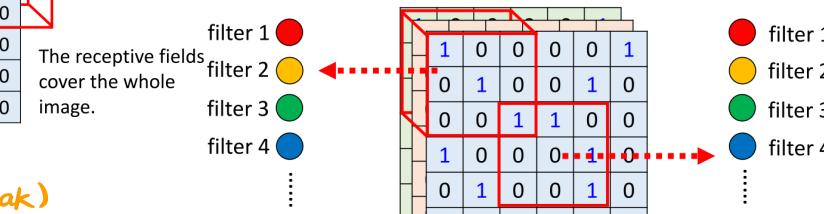
\rightarrow 每个 RF 的 Neuron 参数相同



- Some patterns are much smaller than the whole image.
- The same patterns appear in different regions.

Each receptive field has a set of neurons (e.g., 64 neurons).

Each receptive field has the neurons with the same set of parameters.



Filter 捕捉特征 → Feature Map

- 每个 Filter 检测一个 Pattern
→ 对整张图片做 Inner Product
- Filter 中的 value 是 Parameter

→ Unknown (Gradient Descent 学习行进)

- 一层卷积层会有多个 Filter (e.g. 64)
所有 Filter 的结果 → Feature Map
- Filter 的高度 = Channel 数目

- Deeper Network (More Layers) → 检测到的 Pattern 越大

Feature Map 可以看做 #Channel = #Filter #image

Feature Map 的 Channel 数 = 上一层 image 的 Filter 数

Pooling 池化 通过 Subsampling 减小运算量

- 不是 Layer ← 没有需要学习的参数
更类似 Active Function

- 精度较高的 image 可以取消提高效果

WHY CNN?

- Image 中存在很小的局部特征 (Pattern)

- 相同 Pattern 可能出现在 Image 的不同区域

Subsampling 不会改变物体

Disadvantage: 不能处理 Scaling & Rotation (缩放和旋转) → Data Augmentation

SELF-ATTENTION

应用场景: Sequence as Input

- 文本 One-hot Encoding Vector 间缺乏关联性, 且过于冗长
Word Embedding → 给每个 word 一个向量
一个句子 → 一堆长度不一的向量

- 语音 Frame (语音上的) 向量 e.g. 用一个向量描述 25ms 内的语音信号

Graph 中每个 Node 为一个向量

- Output
- Each Vector has a label POS Tagging 词性标注
 - The whole sequence has a label Sentiment Analysis 情绪判断
 - Model decides #labels Sequence to sequence Translation 翻译

Neighbor Information 上下文信息

全连接层考虑 Neighbor (上下文) 需要很大的 window

→ 容易 Overfitting

Self-Attention 可用作 Input 也可作 Hidden Layer

Structure 结构

- 计算 query, key, value
分别对 a 乘上 W^q, W^k, W^v

- Input vector 间的关联性
 $\alpha \rightarrow$ vector 间的关联性 (关联程度) Attention Score

- 通过 Dot-product 计算 Attention Score α
使用 Softmax 或 ReLU 对其进行 Normalization

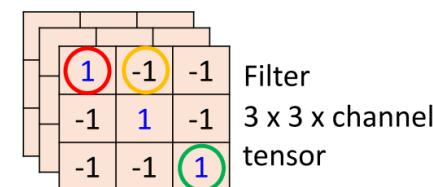
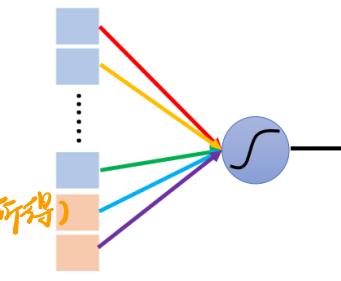
- 基于 Attention Scores 提取信息
 α 越大则 a 对应的 v 会主导 Output vector b

Features 特性

- 从一组 Input Vector (a) 得到一组 Output Vector (b)

- Output vectors ($b_1 \dots b_4$) 同时被计算出来
(Parallel)

- 只有 W^k, W^q, W^v 是需要学习的参数 (未知)



Receptive field

Filter 中有 Bias
(ignore bias in this slide)

Neuron Version Story

Each neuron only considers a receptive field.

Filter Version Story

There are a set of filters detecting small patterns.

The neurons with different receptive fields share the parameters.

Each filter convolves over the input image.

不同的 Neuron 权值共享 观察不同的 Receptive Fields

每个 Filter 对整个 image 做卷积 → Convolution

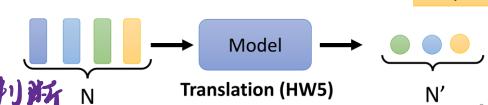
- Each vector has a label. focus of this lecture



- The whole sequence has a label.

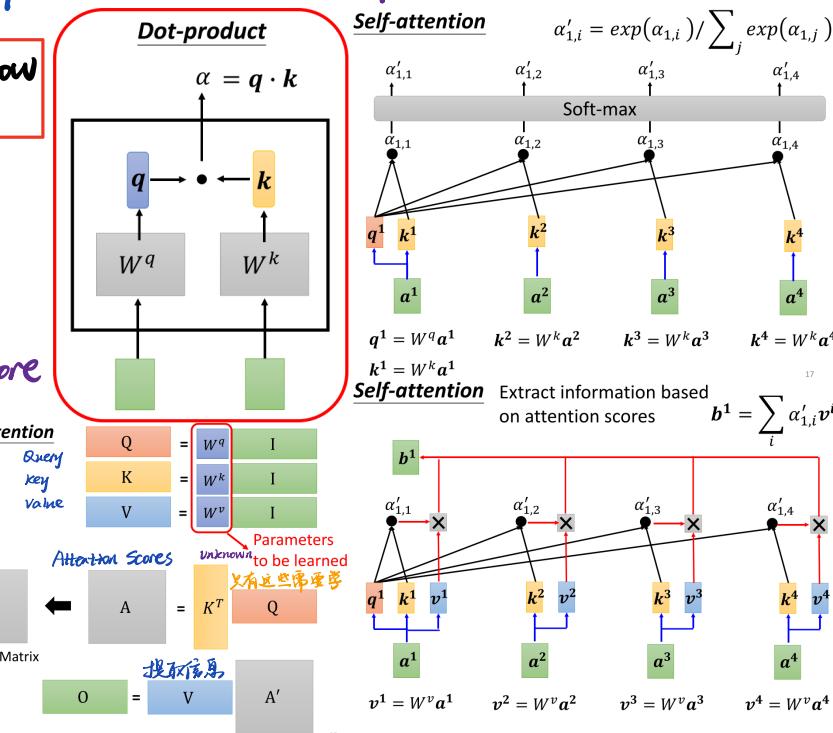


- Model decides the number of labels itself.



seq2seq

Translation (HW5)



Representative Key 去除 Attention Matrix 中的冗余信息

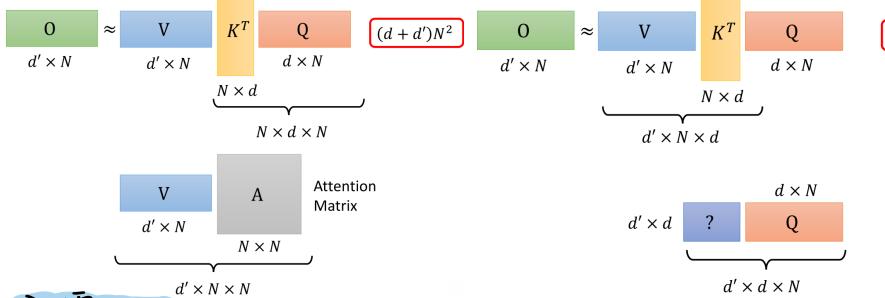
Redundant Columns (Low Rank) \rightarrow 信息冗余

\rightarrow 从 N 个 key 中选取 K 个有代表性的 key & value (key & value 一一对应)
改变 Query 数目 \rightarrow 改变了 Output Sequence 的长度

Compressed Attention 用 CNN 处理 Sequence

Linformer 将 Sequence 看作 $d \times N$ 矩阵, 乘上 $N \times K$ 矩阵 $\rightarrow d \times K$ 矩阵

K, q first $\rightarrow v, k$ first 三个矩阵相乘 \Leftrightarrow Self-Attention



无需 q, k Synthesizer

\rightarrow 将 Attention Matrix 作为网络参数学习 不会使表现变差

\rightarrow Input 不同的 Sequence, Attention Weight 都相同

Attention-free 无需 Attention 处理 Sequence

- Fnet: Mixing tokens with fourier transforms
- Pay Attention to MLPs
- MLP-Mixer: An all-MLP Architecture for Vision

TRANSFORMER

Sequence to Sequence Output Sequence 的长度由模型决定

大部分 NLP 任务可以看作 Question Answering (QA)

\rightarrow 使用 Seq2seq 模型解决 但不一定效果最好

e.g. Google Pixel 4 使用 RNN Transducer

Syntactic Parsing 文法剖析 树状结构 \rightarrow Sequence

e.g. Grammar as a Foreign Language

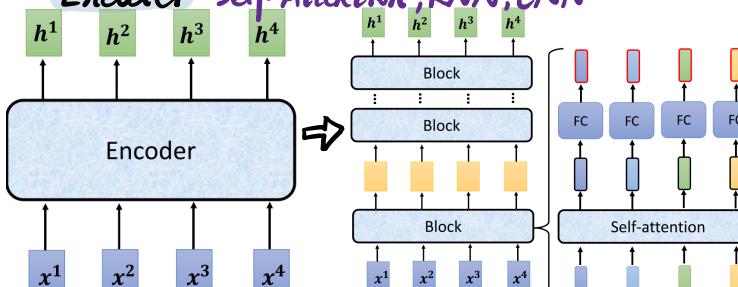
Multi-label Classification 多标签分类

- 同时属于多个 Class
- 由 Seq2seq 模型决定

Object Detection 目标检测 <https://arxiv.org/abs/2005.1287>

Structure 结构

• Encoder self-Attention, RNN, CNN



\rightarrow 输入一个 Vector 序列, Output 一个 Vector Sequence

\rightarrow BERT 和 Transformer 使用了相同的 Encoder

• Decoder Autoregressive (AT)

\rightarrow 读取 Encoder 的 Output 作为 Decoder 的 Input

Bos Begin of Sentence

用 One-hot Vector 表示, 1 是 1, 其他是 0

\rightarrow Decoder 决定 Sequence 的长度

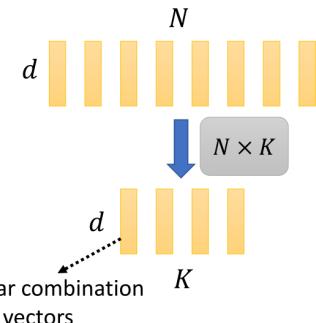
加入 Stop Token END 的概率最大

Compressed Attention

<https://arxiv.org/abs/1801.10198>

Linformer

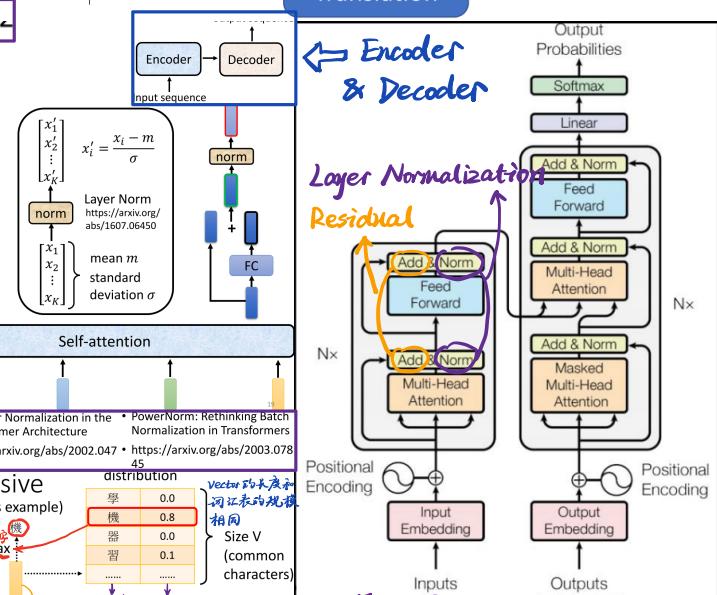
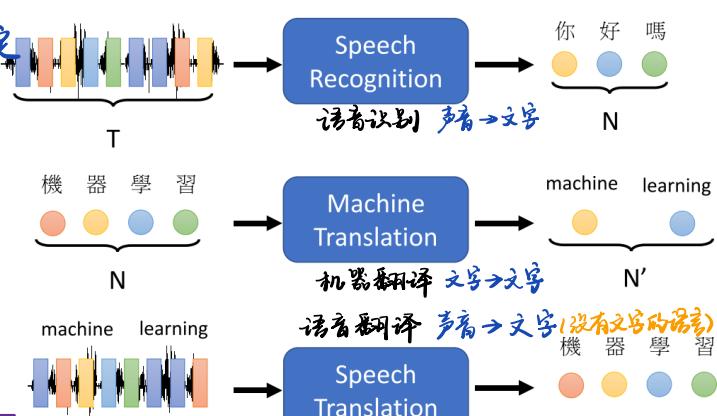
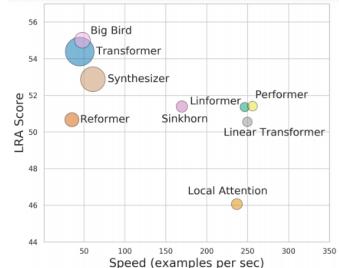
<https://arxiv.org/abs/2006.04768>



CONCLUSIONS

- Human knowledge
 - Local Attention, Big Bird
- Clustering
 - Reformer
- Learnable Pattern
 - Sinkhorn
- Representative key
 - Linformer
- k, q first $\rightarrow v, k$ first
 - Linear Transformer, Performer
- New framework
 - Synthesizer

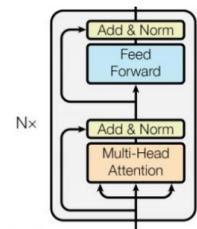
- 图越大, 占用空间越大
- 越右越快
- 越上越准



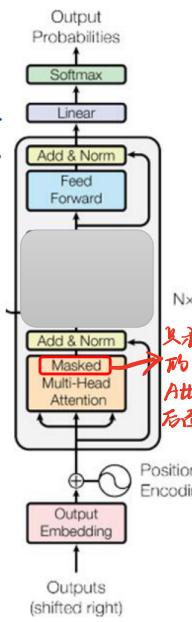
Transformer

<https://arxiv.org/abs/1706.03762>

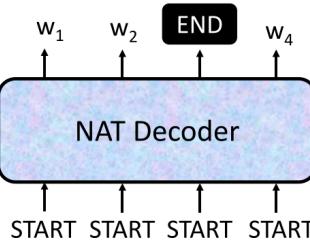
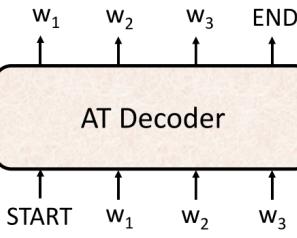
Encoder 的 Input 同时输入
Decoder 的 Input 是一个一个产生的
→ 只能看到当前时刻已有的 Vector
→ Masked Multi-Head Attention



Encoder



Non-autoregressive (NAT) 逐字 Output → Sequence Output
AT v.s. NAT

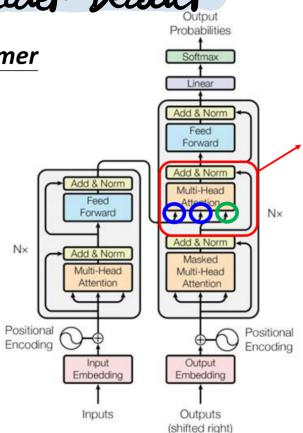


- AT 一个字一个字产生每次一个 Token
- NAT 一词放入一排 BOS, 同时产生整个 Sequence 每次一个 Sequence
- Output Length → 使用另一个 Classifier 输出 Sequence 长度
→ 输出一个长序列，忽略 END Token

Decoder

- Advantage • 并行运算运行更快
- 便于控制 Output Sequence 长度 稳定生成
- NAT 的 Performance 一般差于 AT Multi-modality

Encoder-Decoder Transformer



Training

→ Ground Truth 是一个长度为 vocab 大小的 One-hot Vector

Output 是一个 Probability Distribution

→ 计算每个 Vector 和 GT 的 Cross Entropy

所有 Output 的 Cross-Entropy 越高越小越好

→ 训练时，在 Decoder 的 Input 阶段给 Ground Truth

→ Teacher Forcing 使用 Ground Truth 作为 Input

Tips

Copy Mechanism Create → Copy

有时不需要 Decoder 产生全部的 Output

只需要从 Input 中复制部分作为 Output

→ Pointer Network, Copying Network

Guided Attention 低级错误 (e.g. TTS 漏读字)

引导机器的过程，做 Attention 时设置规则

→ Monotonic Attention, Location-aware Attention

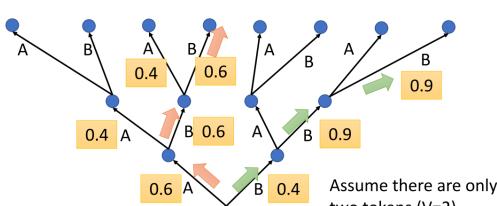
Beam Search 估测出一个最好的 Solution

→ 对有固定答案的任务效果比较好，对需要随机性的任务效果不好

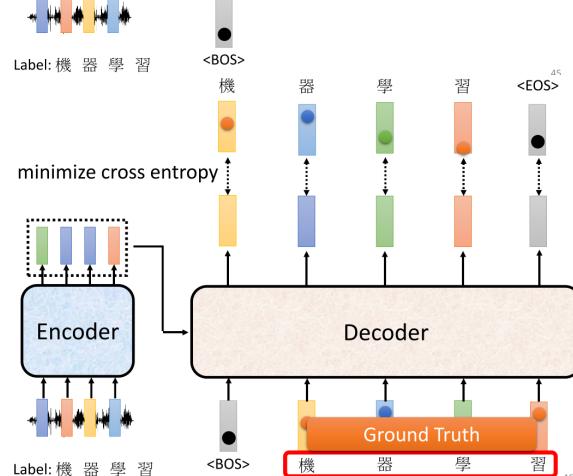
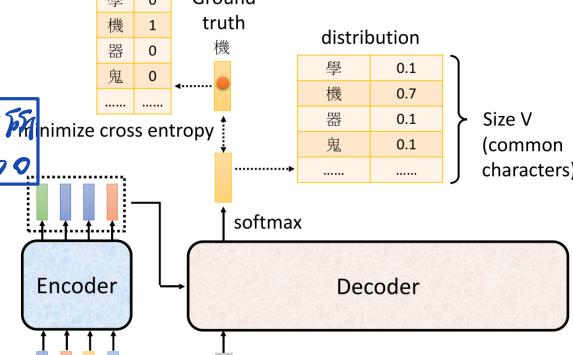
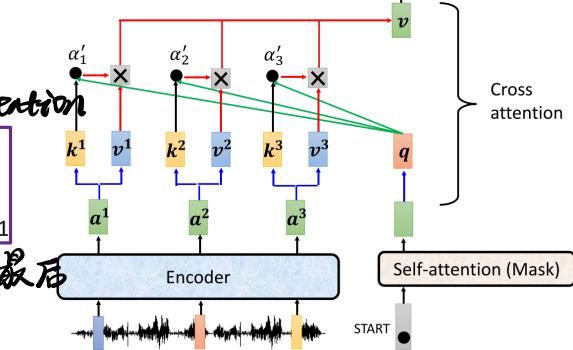
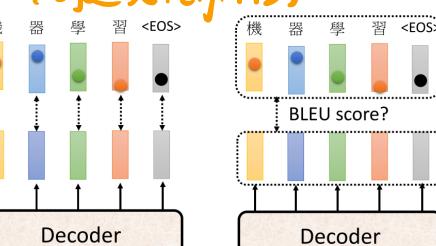
The red path is Greedy Decoding.

The green path is the best one.

Not possible to check all the paths ... → Beam Search



(创造类, e.g. TTS)



Optimization Evaluation Metrics

训练用 Cross-Entropy, 验证用 BLEU

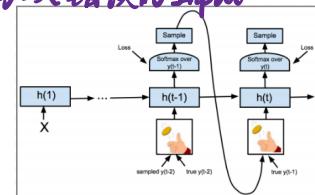
→ 将 BLEU Score 作为 RL Agent 的 Reward
<https://arxiv.org/abs/1511.06732>

Scheduled Sampling

Exposure Bias 训练用的都是 Ground Truth, (Mismatch) Test 时会有错误的 Input 影响准确率
→ 在训练时加入错误的 Input

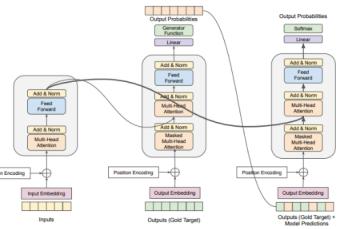
- Original Scheduled Sampling

<https://arxiv.org/abs/1506.03099>



- Scheduled Sampling for Transformer

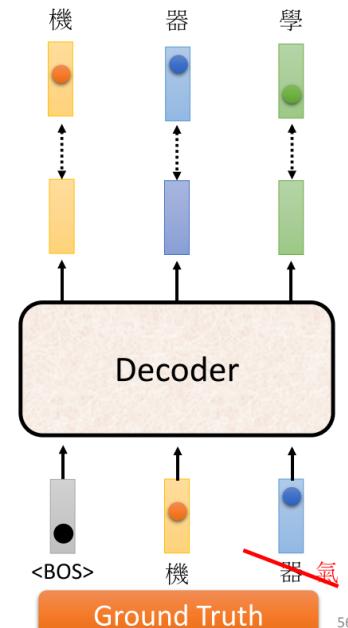
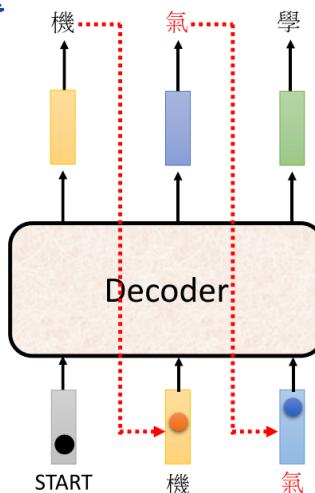
<https://arxiv.org/abs/1906.07651>



- Parallel Scheduled Sampling

<https://arxiv.org/abs/1906.04331>

There is a mismatch! 😕
exposure bias



56

BERT

Self-supervised Learning → Pretrain 产生 BERT 的过程

→ 属于 Unsupervised Learning 的一种 unsupervised 指代消歧方法

- 没有标注 → 取数据的一部分作 Input, 部分作为标注 (x)

- 目标: 令 Output 和标注越接近越好 (y 和 x)

Masking Input <https://arxiv.org/abs/1810.04805> → 学做填空

→ BERT 结构与 Transformer 的 Encoder 相同

- 随机遮住 Sequence 中的一部分 Token → 替换成 special Token
→ 随机换成另一个字

- Output 是一个概率分布 猜测被 Mask 的 Token
分类问题 → 类别与字典规模相同

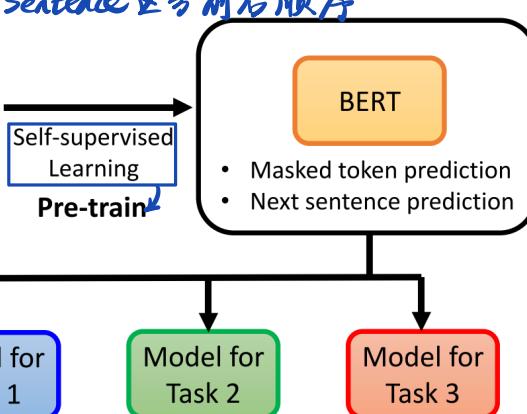
Next Sentence Prediction

- Input 两个句子, 取 CLS 处的 Output 判断两句话是否相接

→ CLS Sequence 开头 SEP 两个句子之间的分隔符

- SOP: Sentence Order Prediction ALBERT

→ 两个 Sentence 区分前后顺序



激发潜能
Fine-tune

How to use BERT

Case 1: Classification

(self)
(supervised)
Pretrain + Finetune semi-supervised

sequence → class

- 取 CLS 对应的 Output 做 Linear 变换后进行分类
- 训练时, Linear 和 BERT 都会进行 Update
- 初始化, Linear 随机, BERT 由 Pre-train 初始化

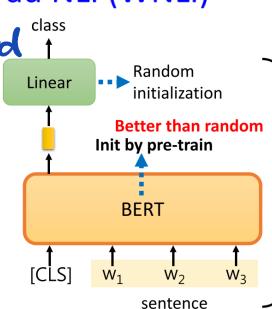
Scratch 整个 Model (BERT & Linear) 都是随机初始化

→ Loss 下降慢且最小值更大

评估 BERT 性能: GLUE General Language Understanding Evaluation

- Corpus of Linguistic Acceptability (CoLA)
- Stanford Sentiment Treebank (SST-2)
- Microsoft Research Paraphrase Corpus (MRPC)
- Quora Question Pairs (QQP)
- Semantic Textual Similarity Benchmark (STS-B)
- Multi-Genre Natural Language Inference (MNLI)
- Question-answering NLI (QNLI)
- Recognizing Textual Entailment (RTE)
- Winograd NLI (WNLI)

Super Glue
更准确
任务集

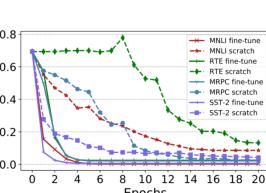


Input: sequence
output: class

Example: 语义分析
Sentiment analysis

this is good
positive

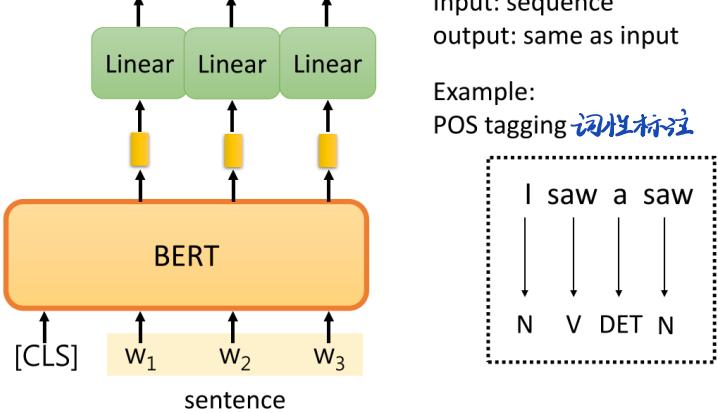
This is the model
to be learned.



实线: Pre-train
虚线: Scratch

Case 2 : Seq2Seq sequence \rightarrow sequence

- Input & Output 都是 Sequence, 且长度相同
- 对 Sequence 的每个 Token 和 Output 做 Classification



Case 4: Question Answering 2 Sequences \rightarrow 2 integers

\rightarrow 给机器一篇文章, 提问 & 回答 Answer - 定在文章中

- Extraction-based Question Answering (QA)

文档: $D = \{d_1, d_2, \dots, d_N\}$

问题: $Q = \{q_1, q_2, \dots, q_M\}$

$$D \rightarrow \text{QA} \rightarrow s$$

$$Q \rightarrow \text{Model} \rightarrow e$$

输出: two integers (s, e)

Answer 第一个字 最后一个字

Answer: $A = \{d_s^1, \dots, d_e^1\}$

In meteorology, precipitation is any product of the condensation of 17 atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain alternate with brief cations called "showers".

What causes precipitation to fall?

gravity $s = 17, e = 17$

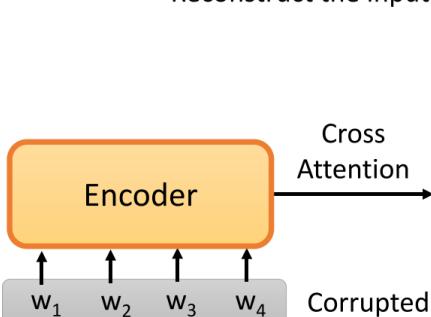
What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

graupel

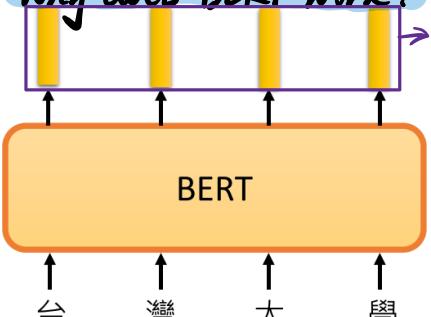
Where do water droplets collide with ice crystals to form precipitation?

within a cloud $s = 77, e = 79$

Reconstruct the input



Why does BERT work?



\rightarrow Embedding Token 的向量化表示

• 意思相近的 token 具有相似的 Embedding

• 考虑上下文 Output 表达 Input 的意思

Contextualized Word Embedding

从上下文可以提取这个 Token 的特征

训练 BERT 微调主题 (Pre-train)

使其具有生成这个位置的表示向量的能力

BERT 用于 DNA 分类 效果非常好 可见 BERT 是一组效果很好的初始化参数

Multi-lingual BERT

English: SQuAD, Chinese: RCDC

Model	Pre-train	Fine-tune	Test	EM	F1
QANet	none	Chinese		66.1	78.1
	Chinese	Chinese		82.0	89.1
BERT	104 languages	Chinese		81.2	88.7
		English		63.3	78.8
		Chinese + English		82.6	90.1

F1 score of Human performance is 93.30%

在多种语言的语料集上训练 BERT

Zero-Shot Reading Comprehension

Pre-train 使用多语言, Fine-tune 使用 English

Test 使用中文, 效果很好

\rightarrow Training: 多语言填空, English 做 QA

Case 3 : 多 Sequence 分类 2 sequences \rightarrow 1 class

判断从前提到否推出假设 Input: two sequences

Output: a class

e.g. 立场分析 Premise Hypothesis

contradiction

entailment

neutral

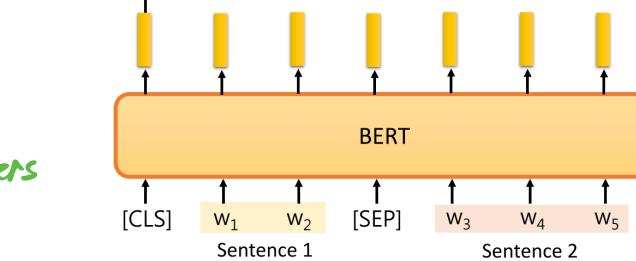
Model

premise: A person on a horse jumps over a broken down airplane

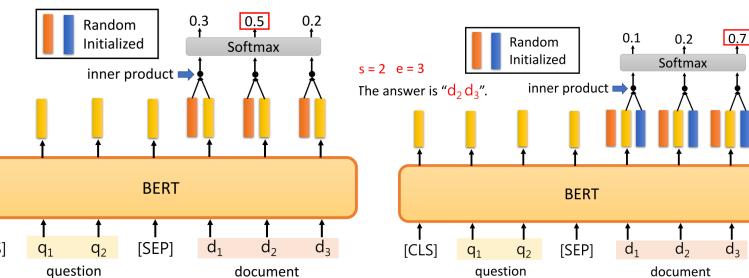
hypothesis: A person is at a diner. contradiction

Class Linear

Input: two sequences Output: a class



- Random Initialize (从头开始训练) 两个人回答
- 预设 Answer 在文档中的开始和结束位置
- Sentence 长度有限制 不超过 512



Pretrain Seq2Seq Model Encoder + Decoder

- Bert & pretrain's Encoder

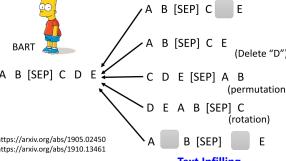
- Pretrain Decoder \rightarrow 破坏 Input 再重建

e.g. Mass . BART

比较各种破坏方法效果: T5

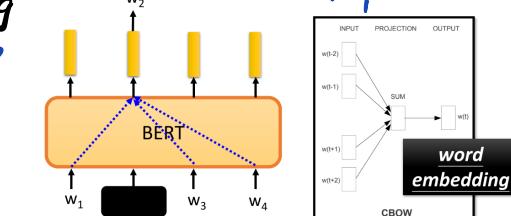
- Transfer Text-to-Text Transformer (T5)
- Colossal Clean Crawled Corpus (C4)

MASS / BART

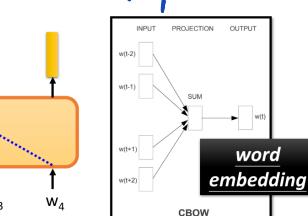


Contextualized word embedding

BERT \leftrightarrow Deep CBOW

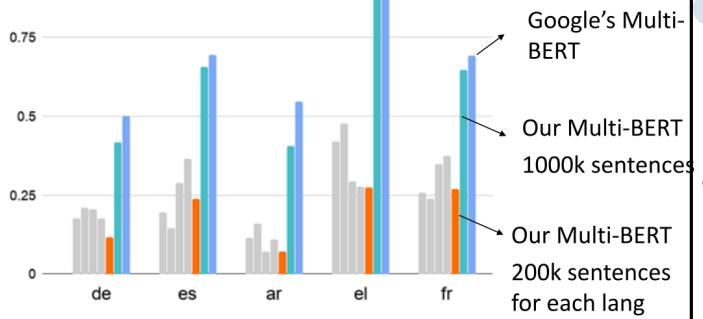


BERT \leftrightarrow Deep CBOW



Cross-lingual Alignment

不同语言的同义词 Embedding 相加



MMR (Mean Reciprocal Rank)

不同语义相同意思的一致程度

• 训练数据规模有较大影响

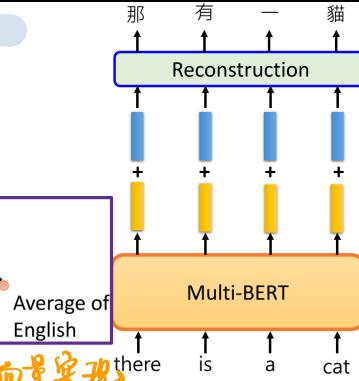
Training Data 越多, MMR 越大

• 能够分辨语言的种类

中文 → 中文 英文 → 英文

• 不同语言可相互转换

→ **Unsupervised Token-level Translation** 通过加入语言间隔向量实现

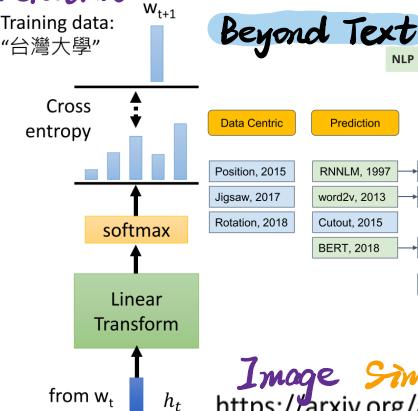
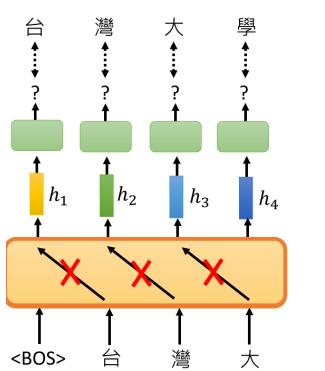


GPT 与 BERT 在 Pre-training 阶段的任务不同

Predict Next Token

- 用所有 Token's Embedding 预测下一个会出现的 Token
- 与 Decoder 相似，但是 Masked 不能看到后面会出现的 Token
- 可以用来生成 Generation

Predict Next Token



How to use GPT?

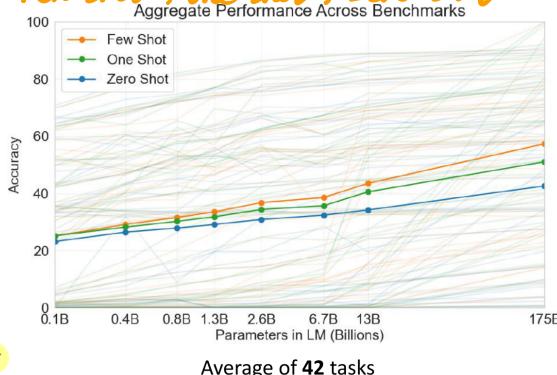
GPT 规模过大, Fine-tune 都较为困难

→ 给出任务的 Description 和少量 example

No Gradient Descent 没有调 GPT 参数

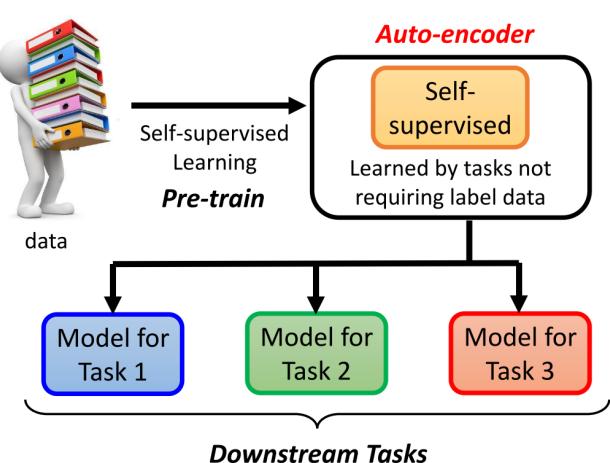
→ "In-context" Learning

Few-shot, One-shot, zero-shot



AUTO-ENCODER

Self-supervised Learning Framework



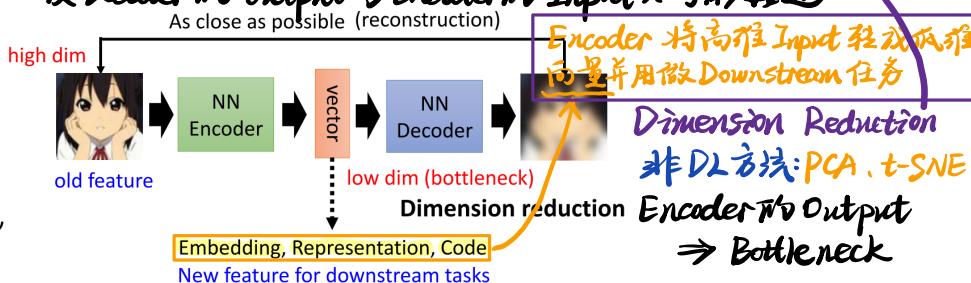
Self-supervised Learning ↔ Pretrain

• 不需要标注 Data 就可以进行训练自行设计训练任务

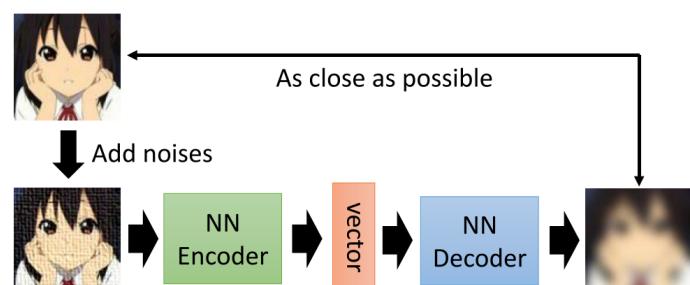
• Auto-encoder 可以看作 Self-supervised Learning 的一种

Reconstruction 重建 用简单方法表示复杂的图片

使 Decoder's Output 与 Encoder's Input 尽可能接近



De-noising Auto-encoder



Feature Disentanglement

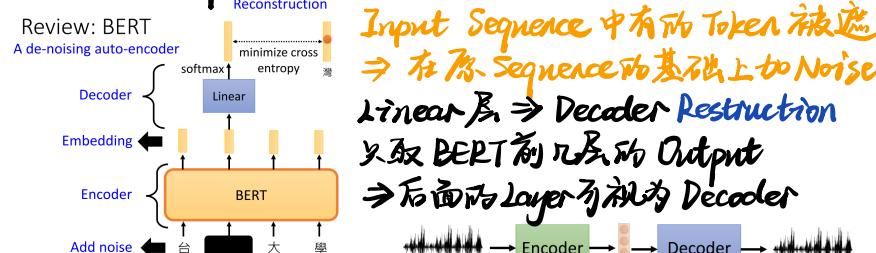
Encoder's Output (Representation) 中包含了 Input 的所有信息

由于高维的 Input 被转成了低维的 Representation 无法分离

Representation 有多维度具体表示了 Input 的哪些信息

→ 在 Auto-encoder 训练时同时知道 Representation 各维度代表了什么信息

给 Encoder's Input 加 Noise, 令 Decoder 抵抗无
Noise 的 Input → 与 BERT 的 mask 机制类似

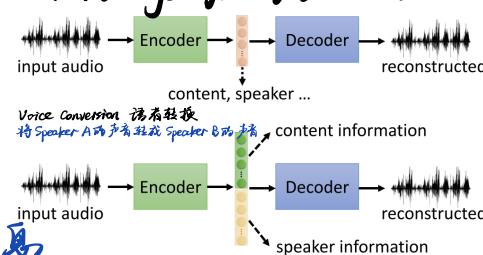


Input Sequence 中有的 Token 被遮
→ 在原 Sequence 的基础上加 Noise

Linear B. → Decoder Reconstruction

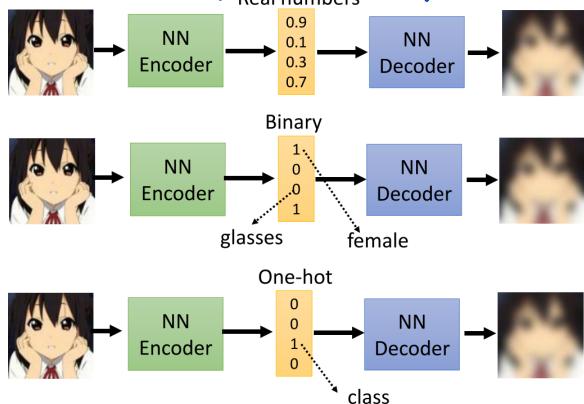
→ 取 BERT 前几层的 Output

→ 后面的 Layer 为 Decoder

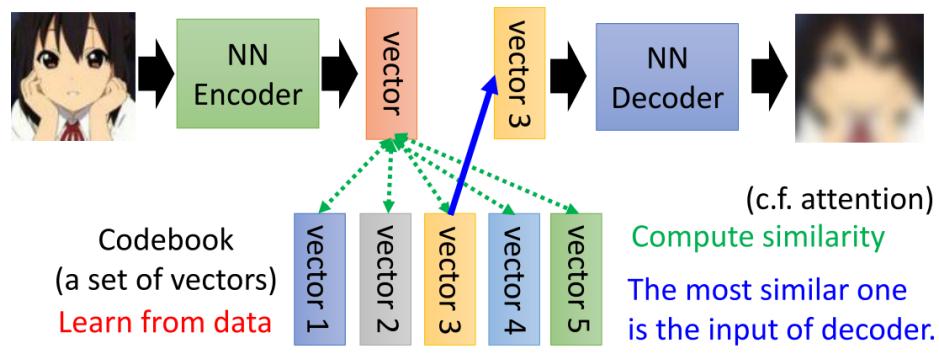


Discrete Latent Representation

Binary 离散向量 Encoder 的输出
One-hot 二进制向量 Encoder 的输出
Unsupervised Classification



• Vector Quantized Variational Auto-encoder (VQVAE)



VQVAE : 元预测的离散化 Codebook 中的 Vector
Phonetic Information <https://arxiv.org/pdf/1901.08810.pdf>

Text as Representation 使用文字代替 Embedding 文章到 Summary

seq2seq2seq auto-encoder → Unsupervised Summarization
中文的 Summary not readable → to Discriminator

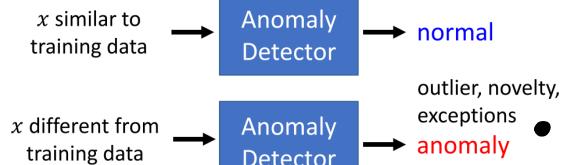
Tree as Embedding <https://arxiv.org/abs/1806.07832>
<https://arxiv.org/abs/1904.03746>

Applications

Generator 从已知 Distribution 随机生成 vector 给 Decoder
产生图片 → VAE (Variational Auto-encoder)

Compression Encoder 压缩, Decoder 扩压 Lossy 图片失真

Anomaly Detection 异常检测 • 检测 Input x 是否与训练数据相似



→ Fraud Detection 信用卡欺诈

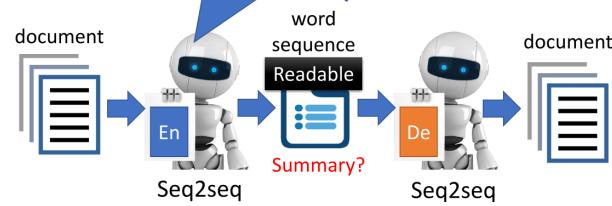
→ Network Intrusion Detection 网络安全

→ Cancer Detection

• One-class Classification Reconstruction Loss 判断 anomaly
绝大部分 Training Data 都是 Positive. 负样本很难收集

Human written summaries → Real or not

Let discriminator considers my output as real



Fraud Detection

- Training data: credit card transactions, x : fraud or not
- Ref: <https://www.kaggle.com/ntnu-testimon/paysim1/home>

Network Intrusion Detection

- Training data: connection, x : attack or not
- Ref: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>

Cancer Detection

- Training data: normal cells, x : cancer or not?
- Ref: <https://www.kaggle.com/uciml/breast-cancer-wisconsin-data/home>