

Contrastive Self-Supervised Learning

Zili Wang

What is Self-Supervised Learning?

Generative / Predictive



Loss measured in the output space

Examples: Colorization, Auto-Encoders

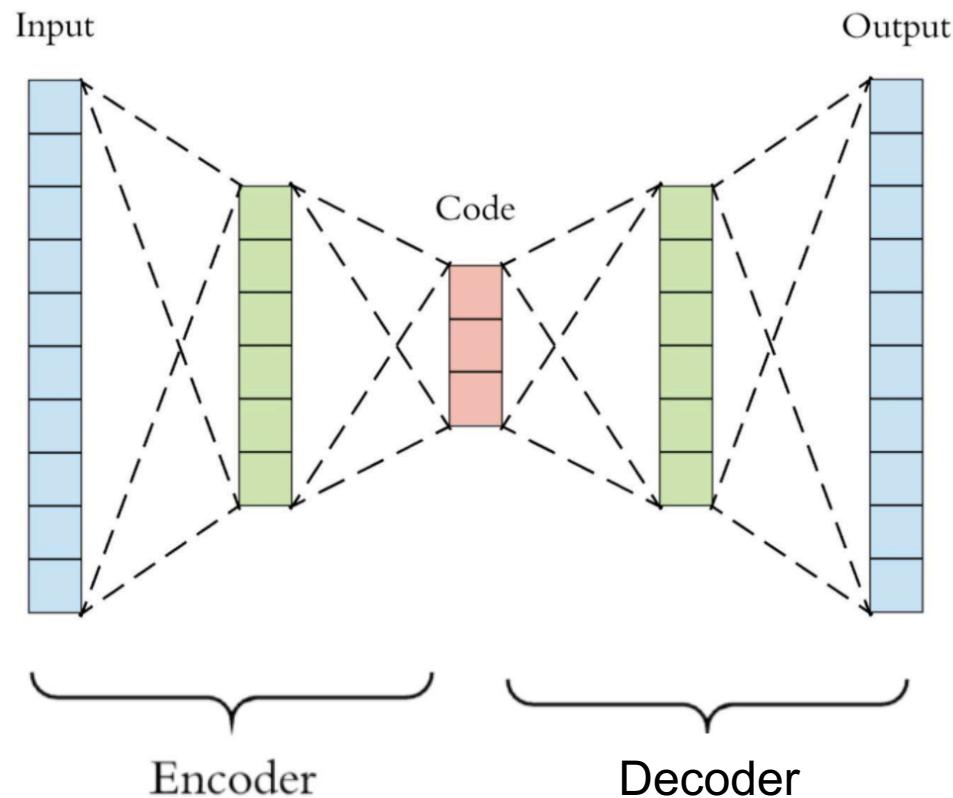
Contrastive



Loss measured in the representation space

Examples: TCN, CPC, Deep-InfoMax

Generative Methods



Contrastive Methods

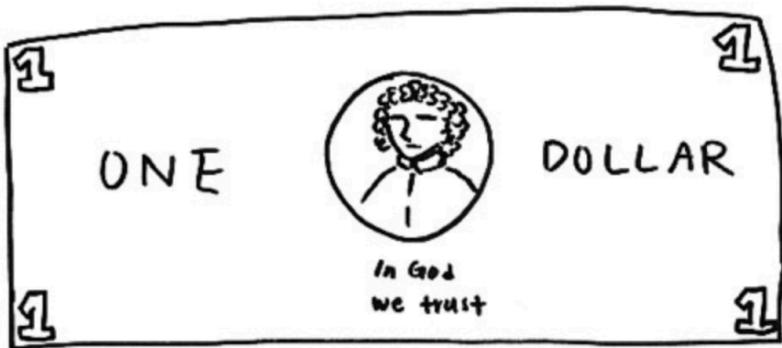


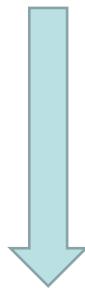
Fig. Left: Drawing of a dollar bill from memory.

Fig. Right: Drawing subsequently made with a dollar bill present.

How do contrastive methods work?

$$\frac{\text{score}(f(x), f(x^+))}{\text{score}(f(x), f(x^-))} >> 1$$

- here x^+ is data point similar or congruent to x , referred to as a *positive* sample.
- x^- is a data point dissimilar to x , referred to as a *negative* sample.
- the score function is a metric that measures the similarity between two features.



Softmax

$$\mathcal{L}_N = -\mathbb{E}_X \left[\log \frac{\exp(f(x)^T f(x^+))}{\exp(f(x)^T f(x^+)) + \sum_{j=1}^{N-1} \exp(f(x)^T f(x_j))} \right]$$

Contrastive methods using margin (2006)

Data: $\mathcal{I} = \{X_1, \dots, X_P\}, X_i \in \mathbb{R}^D$ Goal: G_W

Euclidean distance : $D_W(\vec{X}_1, \vec{X}_2) = \|G_W(\vec{X}_1) - G_W(\vec{X}_2)\|_2$

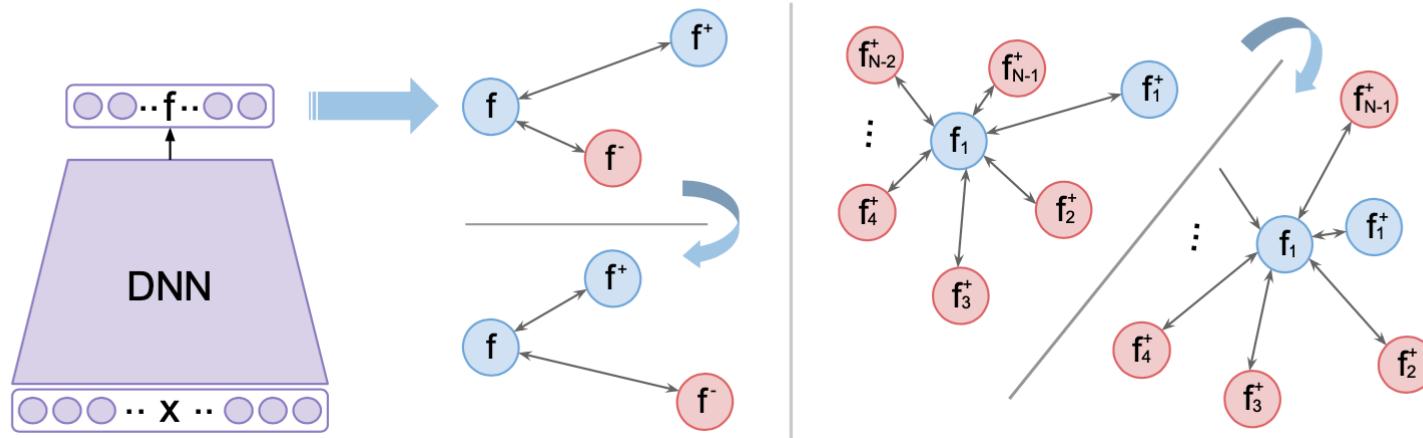
Loss function: $\mathcal{L} = \sum_{i=1}^P (1 - Y) L_S(D_W^i) + Y L_D(D_W^i) ,$

$$L(W, Y, \vec{X}_1, \vec{X}_2) = (1 - Y) \frac{1}{2} (D_W)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_W) \}^2$$

Where the m is the margin

Dimensionality Reduction by Learning an Invariant Mapping
 Raia Hadsell, Sumit Chopra, Yann LeCun (CVPR 2006)

N-pair loss (2016)

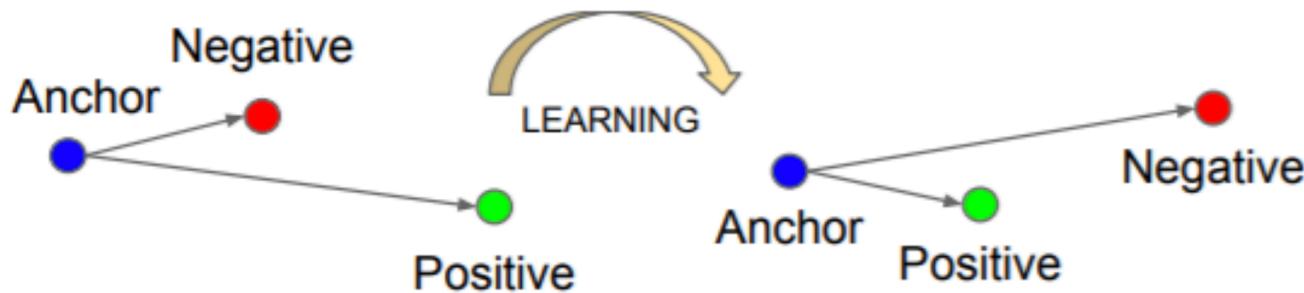


$$\mathcal{L}_{\text{cont}}^m(x_i, x_j; f) = \mathbb{1}\{y_i = y_j\} \|f_i - f_j\|_2^2 + \mathbb{1}\{y_i \neq y_j\} \max(0, m - \|f_i - f_j\|_2)^2$$

$$\mathcal{L}_{\text{tri}}^m(x, x^+, x^-; f) = \max \left(0, \|f - f^+\|_2^2 - \|f - f^-\|_2^2 + m \right)$$

Improved Deep Metric Learning with Multi-class N-pair Loss Objective
 Kihyuk Sohn (NeurIPS 2016)

What is Triplet loss?



$$L = \max(d(a, p) - d(a, n) + \text{margin}, 0)$$

FaceNet: A Unified Embedding for Face Recognition and Clustering
Florian Schroff, Dmitry Kalenichenko, James Philbin citation:5739 2015

N-pair loss (2016)

$$\mathcal{L}_{\text{cont}}^m(x_i, x_j; f) = \mathbb{1}\{y_i = y_j\} \|f_i - f_j\|_2^2 + \mathbb{1}\{y_i \neq y_j\} \max(0, m - \|f_i - f_j\|_2)^2$$

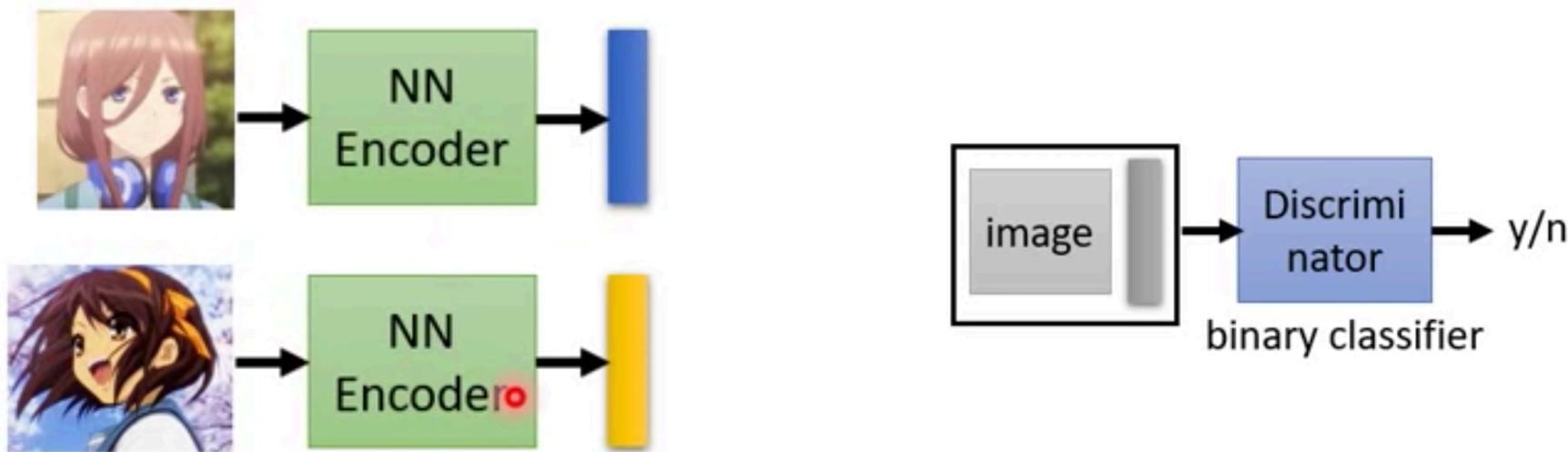
$$\mathcal{L}_{\text{tri}}^m(x, x^+, x^-; f) = \max\left(0, \|f - f^+\|_2^2 - \|f - f^-\|_2^2 + m\right)$$

$$\mathcal{L}_{(2+1)\text{-tuple}}(\{x, x^+, x_i\}; f) = \mathcal{L}_{\text{triplet}}(\{x, x^+, x_i\}; f) = \log(1 + \exp(f^\top f_i - f^\top f^+))$$

$$\log\left(1 + \sum_{i=1}^{L-1} \exp(f^\top f_i - f^\top f^+)\right) = -\log \frac{\exp(f^\top f^+)}{\exp(f^\top f^+) + \sum_{i=1}^{L-1} \exp(f^\top f_i)}$$

Improved Deep Metric Learning with Multi-class N-pair Loss Objective
 Kihyuk Sohn (NeurIPS 2016)

How to evaluate an encoder?



Train ϕ to minimize L_D

$$L_D^* = \min_{\phi} L_D$$

Small L_D^* → The embeddings are representative.

Large L_D^* → Not representative

Contrastive Predictive Coding(2019)

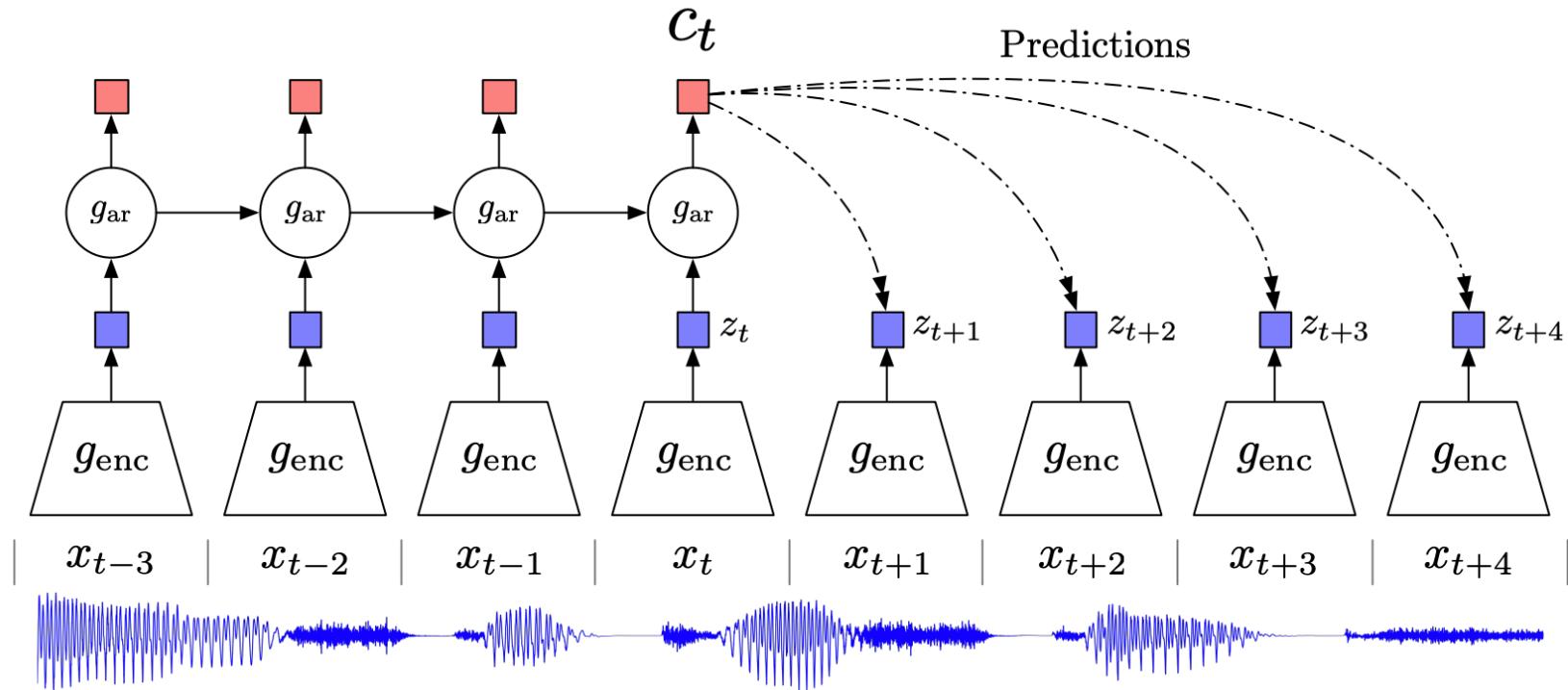


Figure 1: Overview of Contrastive Predictive Coding, the proposed representation learning approach. Although this figure shows audio as input, we use the same setup for images, text and reinforcement learning.

Representation Learning with Contrastive Predictive Coding (Aaron van den Oord, Yazhe Li, Oriol Vinyals) 2018 citation: 385

Contrastive Predictive Coding(2019)

Motivation: Learning Representation

$$z_t = g_{enc}(x_t)$$

$$f_k(x_{t+k}, c_t) \propto \frac{p(x_{t+k} | c_t)}{p(x_{t+k})}$$

Predicating the representation using LSTM

$$f_k(x_{t+k}, c_t) = \exp(z_{t+k}^T W_k c_t)$$

Contrastive Predictive Coding(2019)

Predicating the representation using LSTM

$$f_k(x_{t+k}, c_t) = \exp(z_{t+k}^T W_k c_t)$$

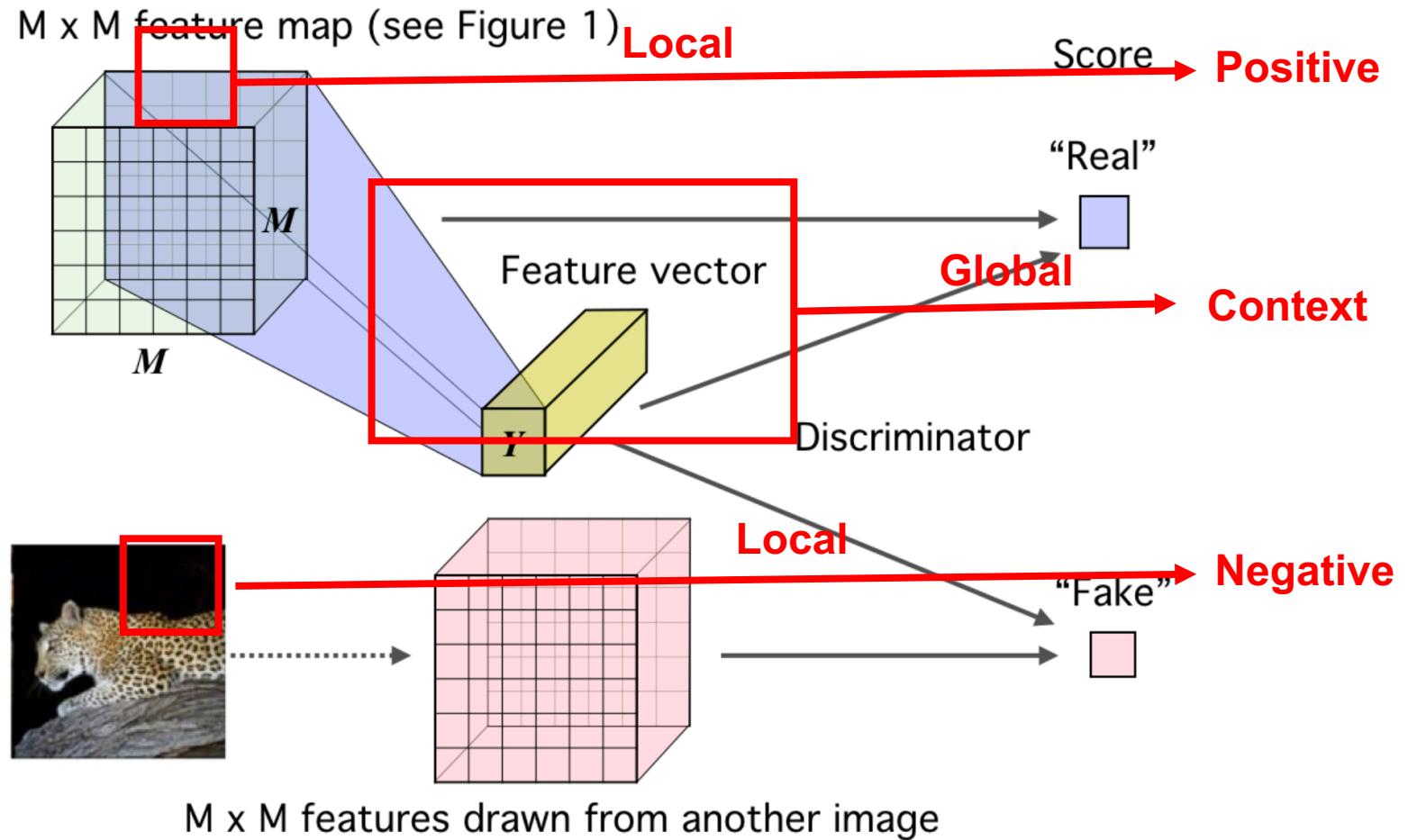
Loss function

$$L_N = -E[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}]$$

$$\frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \quad \longrightarrow \quad \sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

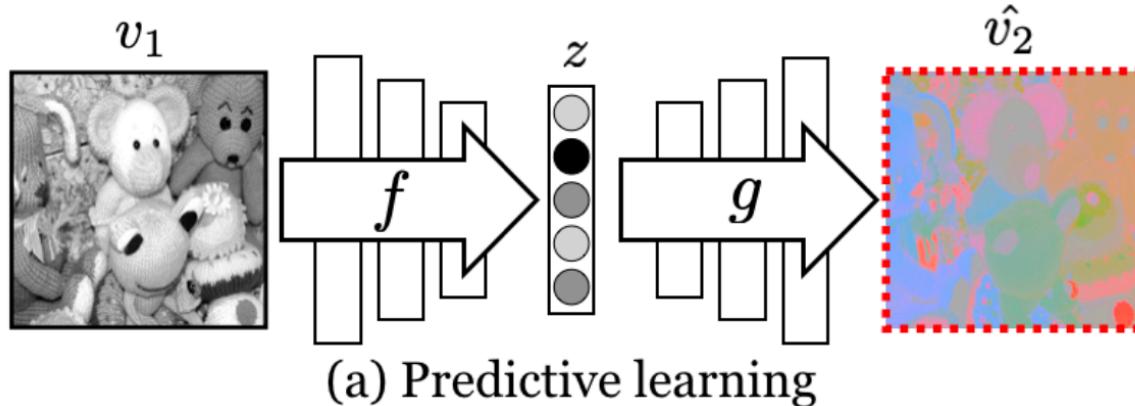
希望这个式子可以接近于1，这样子的话就是loss越来越小，预测的效果越来越好

Learning deep representations by mutual information (2019)



Learning deep representations by mutual information estimation and maximization

Multi-view representation learning (2019)



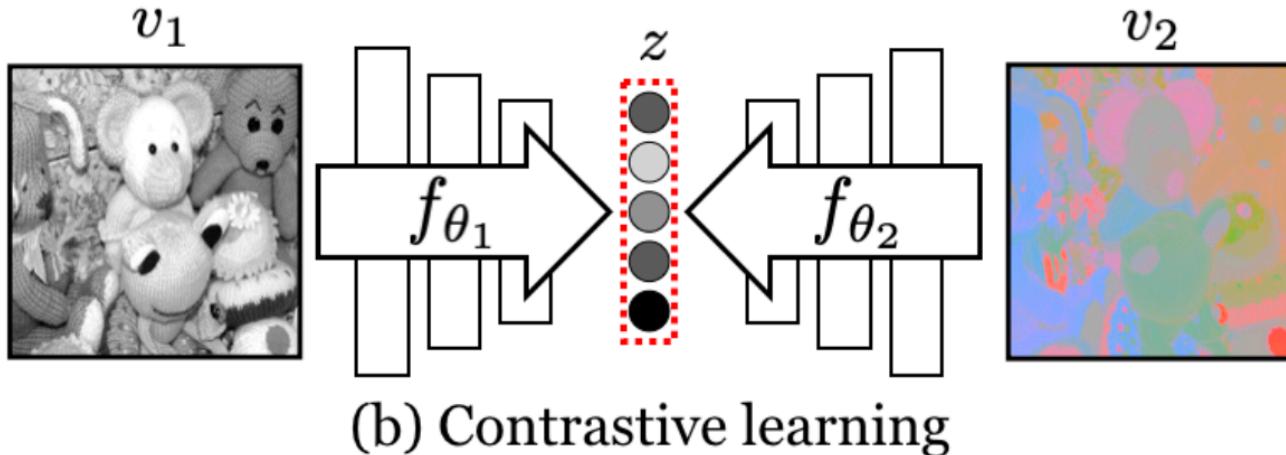
Lab 空间：将图片分解到两个view: L(光度), AB(色度)

v_1 代表光照, v_2 代表色度 $z = f(v_1)$ 和 $\hat{v}_2 = g(z)$

目标：迫使 \hat{v}_2 接近 v_2

缺点：丢失建模关联和复杂结构的能力 \longrightarrow Contrastive Learning

Multi-view representation learning



Positive: $x \sim p(v_1, v_2)$ or $x = \{v_1^i, v_2^i\}$

Negative: $y \sim p(v_1)p(v_2)$ or $y = \{v_1^i, v_2^j\}$

Score function: $h_{\theta}(\cdot)$

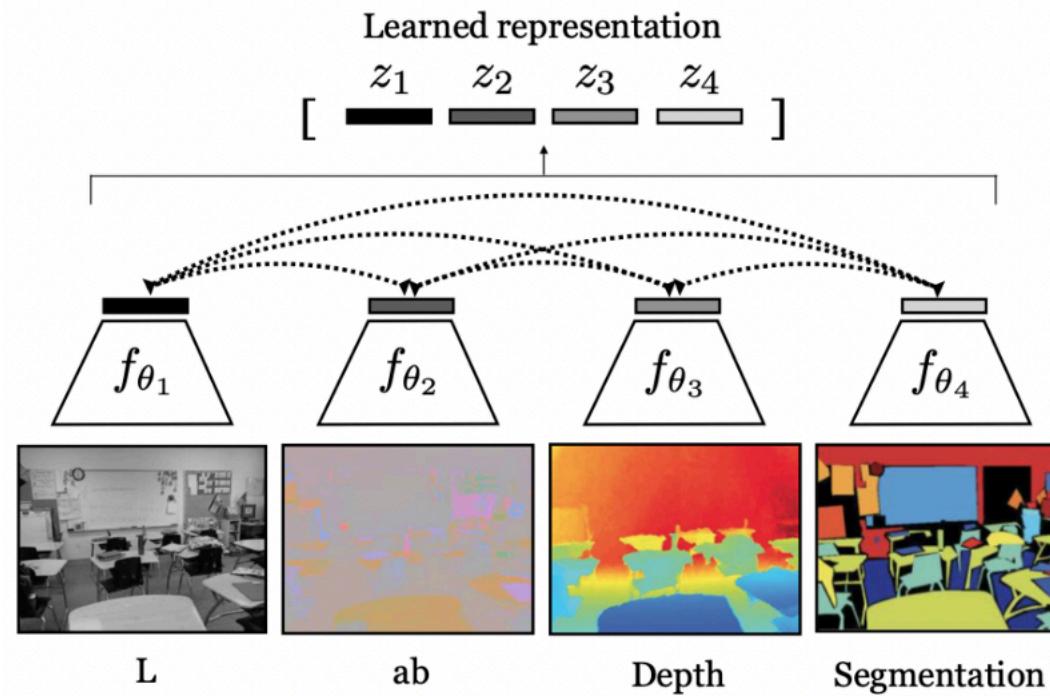
$$\mathcal{L}_{\text{contrast}} = -\mathbb{E}_S \left[\log \frac{h_{\theta}(x)}{h_{\theta}(x) + \sum_{i=1}^k h_{\theta}(y_i)} \right] \xrightarrow{\text{固定一个view}} \mathcal{L}_{\text{contrast}}^{V_1, V_2} = -\mathbb{E}_{\{v_1^i, v_2^i\}} \left[\log \frac{h_{\theta}(\{v_1^i, v_2^i\})}{\sum_{j=1}^N h_{\theta}(\{v_1^i, v_2^j\})} \right]$$

Contrastive Multiview Coding

Yonglong Tian, Dilip Krishnan, Phillip Isola (2019)

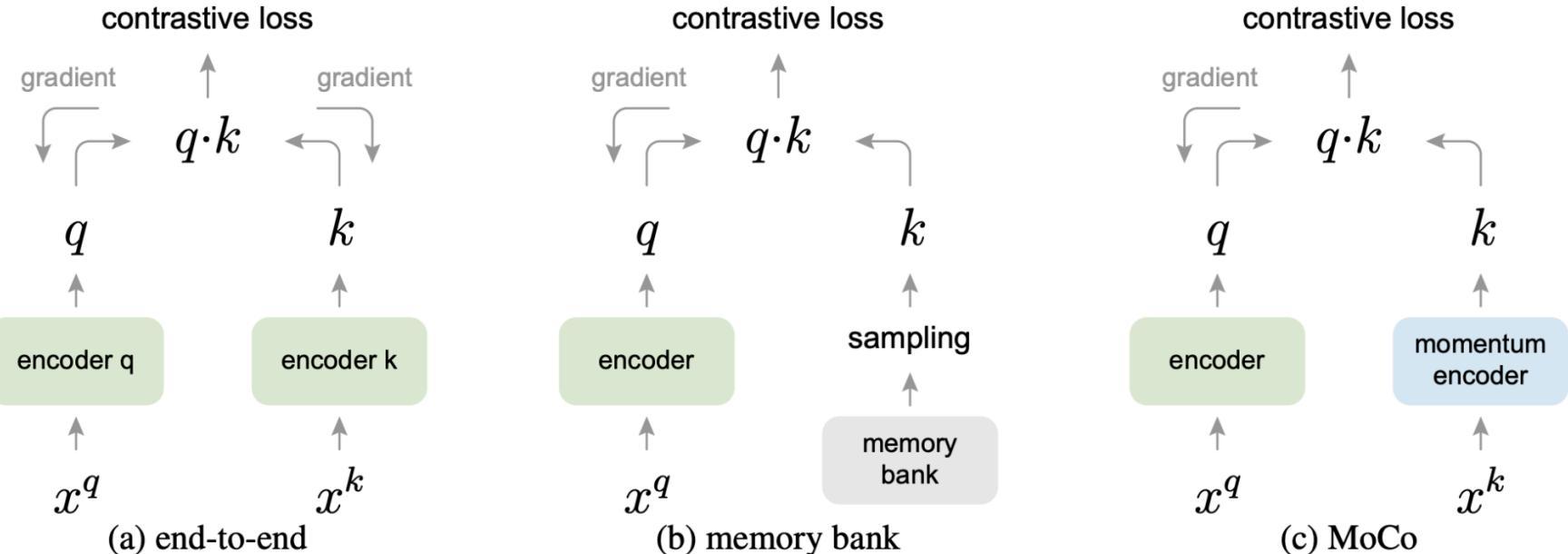
Multi-view representation learning

Learned representations from multiple views are concatenated to form a multi-view learned representation.



 Contrastive Multiview Coding
Yonglong Tian, Dilip Krishnan, Phillip Isola (2019)

Moco



梯度回传
只有当前mini-
batch的数据

1. 需要很大的内存
2. 不能及时更新

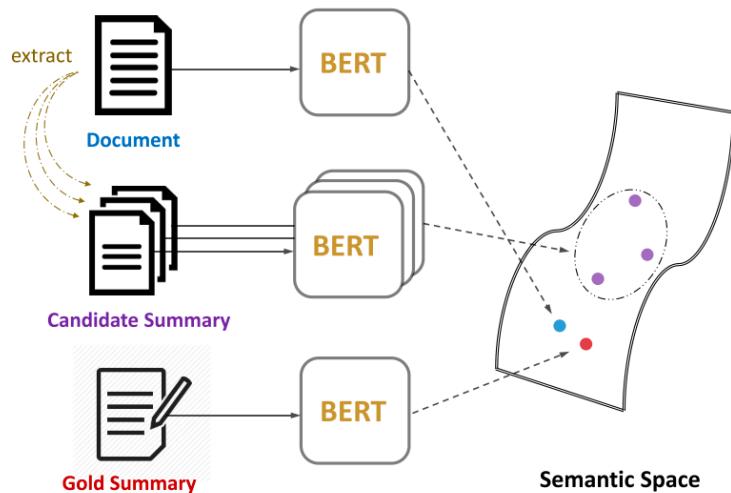
1. 维护一个queue
更新负样本

Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining
Xie Ross Girshick CVPR 2020

Application on Text Summarization

$$\mathcal{L}_1 = \max(0, f(D, C) - f(D, C^*) + \gamma_1), \quad (7)$$



$$\mathcal{L}_2 = \max(0, f(D, C_j) - f(D, C_i) + (j - i) * \gamma_2) \quad (i < j), \quad (8)$$

Model	Multi-News		
	R-1	R-2	R-L
LEAD	43.08	14.27	38.97
ORACLE	49.06	21.54	44.27
MATCH-ORACLE	47.45	17.41	43.14
BERTTEXT	45.80	16.42	41.53
+ 3gram-Blocking	44.94	15.47	40.63
+ 4gram-Blocking	45.86	16.23	41.57
MATCHSUM (BERT-base)	46.20	16.51	41.89

Extractive Summarization as Text Matching
Ming Zhong and Xipeng Qiu et al (ACL 2020)

Application on Machine Translation

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left\{ L(\boldsymbol{\theta}) \right\},$$

$$L(\boldsymbol{\theta}) = \sum_{s=1}^S \log P(\mathbf{y}^{(s)} | \mathbf{x}^{(s)}; \boldsymbol{\theta})$$

$$\hat{\boldsymbol{\theta}}_{\text{CL}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left\{ J(\boldsymbol{\theta}) \right\},$$

where the max-margin loss is defined as

$$J(\boldsymbol{\theta}) = \sum_{s=1}^S \max \left\{ \sum_{n=1}^N \log P(\tilde{\mathbf{y}}_n^{(s)} | \mathbf{x}^{(s)}; \boldsymbol{\theta}) + \eta \right. \\ \left. - N \log P(\mathbf{y}^{(s)} | \mathbf{x}^{(s)}; \boldsymbol{\theta}), 0 \right\}.$$

Reducing Word Omission Errors in Neural Machine Translation: A Contrastive Learning Approach

Zonghan Yang, Yong Cheng, Yang Liu, Maosong Sun(ACL 2019)

Application on Machine Translation

Algorithm 1 Contrastive Learning for NMT

Input: $D = \{\langle \mathbf{x}^{(s)}, \mathbf{y}^{(s)} \rangle\}_{s=1}^S$

Output: $\hat{\boldsymbol{\theta}}_{\text{CL}}$

- 1: Obtain $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ using maximum likelihood estimation on D with random initialization;
 - 2: Construct $\tilde{D} = \{\langle \mathbf{x}^{(s)}, \tilde{\mathbf{y}}^{(s)} \rangle\}_{s=1}^S$ based on D automatically;
 - 3: Obtain $\hat{\boldsymbol{\theta}}_{\text{CL}}$ using contrastive learning on \tilde{D} with $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ as a starting point.
-

$$\hat{\boldsymbol{\theta}}_{\text{CL}} = \operatorname{argmin}_{\boldsymbol{\theta}} \left\{ J(\boldsymbol{\theta}) \right\},$$

where the max-margin loss is defined as

$$J(\boldsymbol{\theta}) = \sum_{s=1}^S \max \left\{ \sum_{n=1}^N \log P(\tilde{\mathbf{y}}_n^{(s)} | \mathbf{x}^{(s)}; \boldsymbol{\theta}) + \eta \right. \\ \left. - N \log P(\mathbf{y}^{(s)} | \mathbf{x}^{(s)}; \boldsymbol{\theta}), 0 \right\}.$$

For each ground-truth sentence pair $\langle \mathbf{x}^{(s)}, \mathbf{y}^{(s)} \rangle$, it is possible to sample N negative examples $\langle \mathbf{x}^{(s)}, \tilde{\mathbf{y}}_1^{(s)} \rangle, \dots, \langle \mathbf{x}^{(s)}, \tilde{\mathbf{y}}_N^{(s)} \rangle$. For simplicity, we set $N = 1$ and use $\tilde{D} = \{\langle \mathbf{x}^{(s)}, \tilde{\mathbf{y}}^{(s)} \rangle\}_{s=1}^S$ in our experiments.

Reducing Word Omission Errors in Neural Machine Translation: A Contrastive Learning Approach

Zonghan Yang, Yong Cheng, Yang Liu, Maosong Sun(ACL 2019)

Application on Machine Translation

Dataset: WMT2017

Systems	EN-DE			ZH-EN	
	BLEU	#Speed	#Param	BLEU	#Param
<i>Existing NMT systems</i>					
Trans.base (Vaswani et al., 2017)	27.3	N/A	65.0M	N/A	N/A
+Context-Aware SANs (Yang et al., 2019a)	28.26	N/A	106.9M	24.67	126.8M
+Convolutional SANs (Yang et al., 2019b)	28.18	N/A	88.0M	24.80	N/A
+BIARN (Hao et al., 2019)	28.21	N/A	97.4M	24.70	107.3M
Trans.big (Vaswani et al., 2017)	28.4	N/A	213.0M	N/A	N/A
+Context-Aware SANs (Yang et al., 2019a)	28.89	N/A	339.6M	24.56	379.4M
+Convolutional SANs (Yang et al., 2019b)	28.74	N/A	339.6M	25.01	N/A
+BIARN (Hao et al., 2019)	28.98	N/A	333.5M	25.10	373.3M
<i>Our NMT systems</i>					
Trans.base	27.48	13.2K	66.5M	24.28	74.7M
+SCWAContext	28.28+	12.1K	72.8M	24.79+	81.0M
+TCWALoss	27.94+	14.3K	66.5M	24.65	74.7M
+BCWAContLoss	28.51+	13.1K	72.8M	24.94+	81.0M
Trans.big	28.45	11.2K	221.1M	24.55	237.5M
+BCWAContLoss	29.14+	10.1K	246.3M	25.12+	262.7M

Fig. [2]

Method	Zh-En	De-En	Ru-En
MLE	23.90	34.88	31.24
MLE + CP	24.04	34.93	31.36
WordDropout	23.73	34.63	31.05
CL _{one}	24.92 ++**††	35.74 ++**††	32.04 ++**††
CL _{two}	24.76 ++**††	35.54 ++**††	31.94 ++*††
CL _{three}	24.52 +*††	35.44 ++*††	32.20 ++**††
CL _{low}	24.13 †	34.96 †	31.47 ++†
CL _{high}	24.77 ++**††	35.24 ++††	31.70 ++††
CL _V	24.12 †	35.02 ††	31.73 ++*††
CL _{IN}	24.71 ++**††	35.26 +*††	31.76 ++*††

Fig. [1]

[1] Reducing Word Omission Errors in Neural Machine Translation:
A Contrastive Learning Approach

Zonghan Yang, Yong Cheng, Yang Liu, Maosong Sun(ACL 2019)

[2] Content Word Aware Neural Machine Translation

 Kehai Chen, Rui Wang*, Masao Utiyama, and Eiichiro Sumita(ACL 2020)

Application on Commonsense Reasoning

Task: Pronoun Disambiguation and Winograd Schema Challenge problems.

Sentence-1: *The trophy doesn't fit in the suitcase because **it** is too small.*

Answers: A) the trophy B) the suitcase

Sentence-2: *The trophy doesn't fit in the suitcase because **it** is too big.*

Answers: A) the trophy B) the suitcase

Application on Commonsense Reasoning

PDP-60 (sup.) (Davis et al., 2016)	
Patric Dhondt (WS Challenge 2016)	45.0 %
Nicos Issak (WS Challenge 2016)	48.3 %
Quan Liu (WS Challenge 2016-winner)	58.3 %
USSM + Supervised DeepNet	53.3 %
USSM + Supervised DeepNet + 3 KB	66.7 %
BERT-ft (Kocijan et al., 2019)	78.3 %
PDP-60 (unsupervised)	
Unsupervised Sem. Similarity (USSM)	55.0 %
Transformer LM (Vaswani et al., 2017)	58.3 %
BERT LM (Trinh and Le, 2018)	60.0 %
MAS (Klein and Nabi, 2019)	68.3 %
DSSM (Wang et al., 2019)	75.0 %
CSS (Proposed Method)	90.0 %

WSC-273 (sup.) (Levesque et al., 2012)	
USSM + KB	52.0%
USSM + Supervised DeepNet + KB	52.8 %
Transformer (Vaswani et al., 2017)	54.1 %
Know. Hunter (Emami et al., 2018)	57.1 %
GPT-ft (Kocijan et al., 2019)	67.4 %
BERT-ft (Kocijan et al., 2019)	71.4 %
WSC-273 (unsupervised)	
Single LMs (Trinh and Le, 2018)	54.5 %
MAS (Klein and Nabi, 2019)	60.3 %
DSSM (Wang et al., 2019)	63.0 %
Ensemble LMs (Trinh and Le, 2018)	63.8 %
CSS (Proposed Method)	69.6 %
DPR (Rahman and Ng, 2012)	
(Rahman and Ng, 2012)	73.0%
(Peng et al., 2015)	76.4 %
CSS (Proposed Method)	80.1 %
KnowRef (Emami et al., 2019)	
E2E (Emami et al., 2019)	58.0 %
BERT-ft (Emami et al., 2019)	61.0 %
CSS (Proposed Method)	65.5 %

Table 1: Results on different tasks. From Top to bottom: PDP, WSC, DPR, KnowRef. The first two task performances are subdivided into two parts. Upper part: supervised, lower part: unsupervised.

Contrastive Self-Supervised Learning for
Tassilo Klein, Moin Nabi(ACL 2020)

Other reference

[1] CONTRASTIVE LEARNING OF STRUCTURED WORLD MODELS

Thomas Kipf, Elise van der Pol, Max Welling (ICLR 2020)

[2] DEEP GRAPH INFOMAX

Petar Veličković, William Fedus, William L. Hamilton, Pietro Liò,
Yoshua Bengio, R Devon Hjelm (ICLR 2019)

[3] Adversarial Contrastive Estimation.

Avishek Joey Bose, Huan Ling, Yanshuai Cao (ACL 2018)

[4] Contrastive Language Adaptation for Cross-Lingual Stance Detection

Mitra Mohtarami, James R. Glass, Preslav Nakov (EMNLP 2019)