

I-JEPA reading notes

1. Introduction

image self-supervised learning:

① invariance-based method:

similar embeddings for ≥ 2 views of two same image

— high semantic level representation ✓

— hand-crafted data augmentation ✗

— strong biases ✗

② generative method:

remove/corrupt portions of the input + predict corrupted content = pixel / token

— less prior knowledge ✓

— generalize beyond modality ✓

— lower semantic level ✗

— underperform ① ✗

2. Methods: I-JEPA framework

<1> targets

input image $y \longrightarrow N$ non-overlapping patches



target encoder $f_{\bar{\theta}}$



patch level representation $s_y = \{s_{y1}, \dots, s_{yN}\}$

↓ random sample

M blocks: $M=4$

aspect ratio = (0.75, 1.5), scale = (0.15, 0.2)

the i th block cover mask $B_i: s_y(i) = \{s_{yj}\}_{j \in B}$

representation

0 or 1

output of target encoder

<2> Context

input image \rightarrow single block x + mask B_x :
unit aspect ratio + random scale (0.85, 1.0)
 \rightarrow remove overlap with M targets

\downarrow
context encoder f_θ

\downarrow
patch level representation $S_x = \{s_{x_j}\}_{j \in B_x}$

<3> Prediction

input: $S_x + \{m_j\}_{j \in B_i}$ \nearrow conditioned on
a mask token for each patch we wish to predict

\downarrow
 g_ϕ

output: $\hat{S}_y(i) = \{\hat{s}_{y_j}\}_{j \in B_i} = g_\phi(S_x, \{m_j\}_{j \in B_i})$
apply predictor M times: $\hat{S}_y(1), \dots, \hat{S}_y(M)$

* mask tokens: parameterized by a shared learnable vector
with an added positional embedding

* Loss: L_2 loss between \hat{S}_{y_j} and S_{y_j} .

* ϕ, θ learned through gradient descent

$\bar{\theta}$ updated via exponential moving average of θ

3. Self-supervised learning architectures

objective: Energy-Based Models (EBM)

incompatible inputs \longrightarrow high energy

compatible inputs \longrightarrow low energy

<1> Joint-Embedding Architectures (JEA)

invariance-based pretraining

incompatible inputs \longrightarrow dis-similar embeddings

compatible inputs \longrightarrow similar embeddings

Challenge: representation collapse

the encoder produces a constant output regardless of the input

<2> Generative Architectures

reconstruction-based methods

produce compatible x, y conditioned on z position tokens

\downarrow
copy of y with mask

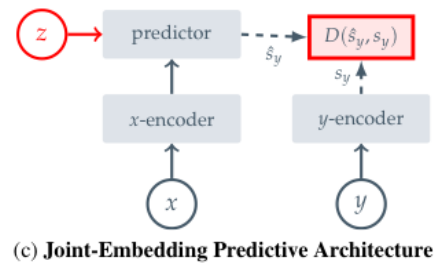
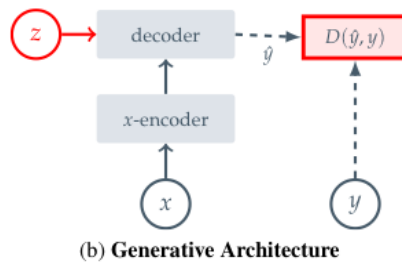
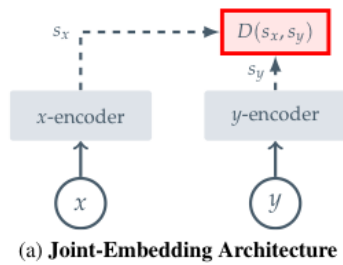
No representation collapse: the informational capacity of z is low compared to the signal y

<3> Joint-Embedding Predictive Architecture (JEPA)

predict embedding of y from a compatible signal x
conditioned on z

Challenge: representation collapse

solution: asymmetric architecture between x - and y -encoders.



4. Other

<1> I-JEPA有强语义性的2个原因:

- ① unnecessary pixel level details are potentially eliminated.
- ② multi-block masking strategy

<2> 接近的工作: data2vec

Context Autoencoders

AI-JEPA

noise I-JEPA