

## Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
  - A. 이름: **DenoiseRep: Denoising Model for Representation Learning**
  - B. 저널: **NeurIPS**
  - C. 도메인: **Representative**
  - D. 출판연도: **2024**
  - E. 저자: **Zhengrui Xu, Guan'an Wang, Xiaowen Huang, Jitao Sang**
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. 논문 **Figure**를 그대로 따라 그리면 안됩니다.
  - A. **Feature**을 잘 추출하는 것은 **Discriminative Model**에게 있어서 매우 중요한 업무라고 할 수 있음.
  - B. 이를 효과적으로 추출하기 위해 우수한 성능의 모델인 Denoise Model의 구조 (Diffusion)를 적용하고자 함.
  - C. 일반적으로 사용되는 Feature Extract backbone (예: Res50, ViT, CNN)에 우수한 성능을 보인 Diffusion과 같이 Noise를 주입하고 이를 제거하는 과정을 거침으로써, Robust 한 Feature를 추출할 수 있다는 판단을 함. (가설)
  - D. 그러나, 매번 추출 후, Forward – Reverse과정을 거치는 과정이 필요하고 이는 모델을 여러 번 호출하게 만드는 Latency문제를 야기함.
  - E. 따라서, Embedding (Backbone)과정과 Denoising (Diffusion) 과정을 한 step 만에 해결할 수 있는 아키텍처를 제안함. (이때, 연산량 + Parameter의 증가를 해결할 수 있는 구조를 제안)
  - F. 일단, Backbone에 해당하는 부분의 경우 Train하지 않고 본인의 업무를 하도록 둠.
  - G. Denoise 부분에 대해서만 Train하도록 함 (L2 Loss, Embed Layer의 Output을 Input으로 사용하며, Input에 대해 Forward Process – Reverse Process, 이

때 Noise를 Predict하는 식으로 학습. 해당 구조는 Linear한 구조로, 한 번에 1개의 Denoise step이 아닌, 2-n개씩 처리가능한 확장이 가능한 구조임. 이 때, Noise는 20회 정도 주입한다고 되어있음.)

3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
  - A. 해당 부분을 적용하는데 난이도가 상당히 쉽다는 것이 가장 큰 장점
  - B. 앞서 보인 아키텍처를 통해 얻은 Robust Feature을 통해 같은 Task를 진행했을 때 성능의 향상을 보였다는 점은 고무적임.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
  - A. 이 논문을 계기로 Representation Generation Method의 다양한 논문을 Review하게 될 예정임.
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
  - A. 다만, Diffusion을 feature에도 적용한다? 연산이 추가적으로 필요한 구조라는 단점을 가지고 있음
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
  - A. 다양한 Modality를 사용하는 상황에서 이를 Ablation 실험에 시도하는 것은 좋다고 생각됨. 그러나 더 많이 고려해볼 필요가 있다고 생각함.
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
  - A. 아니요

날짜: 2025-08-17

이름: 신준원

# DenoiseRep: Denoising Model for Representation Learning

## Introduction

denoising model  $\rightarrow$  discriminative task을 적용해볼까.

$\rightarrow$  Representation learning 이 필요함.

$\Rightarrow$  왜 중요한 것일까?

$\Rightarrow$  data의 representation을 잘 학습한 경우  $\rightarrow$  정보 추출이 용이할 수 있음.

$\rightarrow$  Classification or other predictor 등

여기서, 불완전한 denoising을 활용한 효과적인 Representation task 활용 방법을 어떻게 할까?

$\rightarrow$  feature discrimination

$\oplus$  joint feature extraction & denoising

## Generative model (with denoising)

DDPM :  $P(x, y)$   
 $\downarrow$   
Sample Condition

## Generative vs Discriminative

denoising 적용  $\Rightarrow$  Generative model > Discriminative model

discriminative model :  $P(y|x)$

$\downarrow$   
label, etc

$\rightarrow$  이, 같은 diffusion model :

- ① DiffusionDet : Diffusion + Object Detection
- ② DiffSeg : Diffusion + Image Segmentation

이런  $\rightarrow$  ex) noise box, noise segmentation

$\Rightarrow$  learning representation  $\Rightarrow$  위한 정보 추출에 용이함

$\rightarrow$  task에 활용 가능 (with denoising model)

task  $\Rightarrow$  ReID  $\rightarrow$  이 feature 추출함.

$\rightarrow$  cloud 등?

## Denoise Rep

→ Representation learning. → 노이즈를 제거하는 방법

① discriminative model's backbone (Vision) : ResNet, ViT

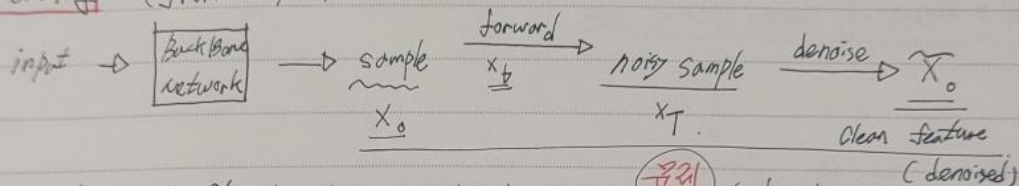
→ feature extraction

② specific head : MLP for classification  
RCNN for object detection  
FCN for segmentation

→ Denoising layer → Introduced

→ 노이즈를 제거한 feature extraction

## Denoising (from DDPM)



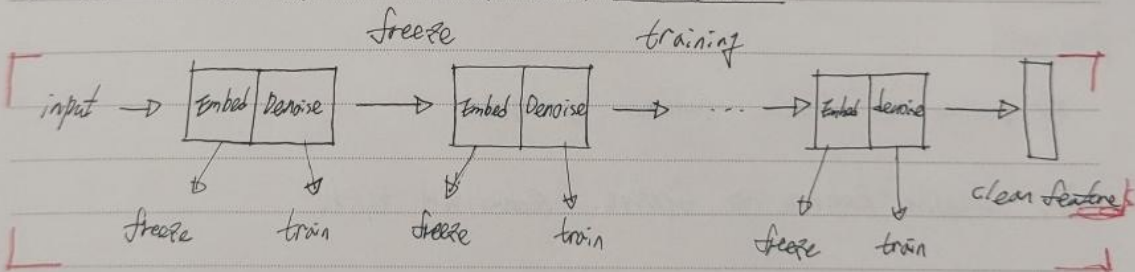
⇒  $X_0$  대비  $\tilde{X}_0$  의 특징이 제거된 특징을 추출. (노이즈 제거)

Denoise T step

⇒ denoise step 필요

→ Backbone + denoising layer fusion

⇒ 2 branch (original embedding layer, denoising layer)



input → task 1 (denoise) → fusion → task 2, task 3

## Feature extraction & Denoising

input  $x$ 을 Diffusion Network 이용  $x$ .  
→ feature (from backbone) 를 input으로 사용하는 것.  $\frac{x_0}{4}$   $\rightarrow$  1st input  $\rightarrow$  2nd input.  
feature의 Denoising  $\Rightarrow$  robust feature를 위한 것.

예를 들어 feature  $\rightarrow$  output  $\rightarrow$  noise  $\rightarrow$  denoising (recursive!)  
Backbone  $\rightarrow$  Diffusion  $\rightarrow$  Denoising layer  $\rightarrow$  feature.

Step: 이름을 inference step.

$\rightarrow$  feature  $\rightarrow$  diffusion model  $\rightarrow$  denoising layer into next embedded layer.  
 $\rightarrow y = Wx + b$ .

## Experiments

only backbone VS Denoise backbone  $\Rightarrow$  성능

즉, Denoise backbone  $\Rightarrow$  feature (Robust) 하기 성능