

MIL

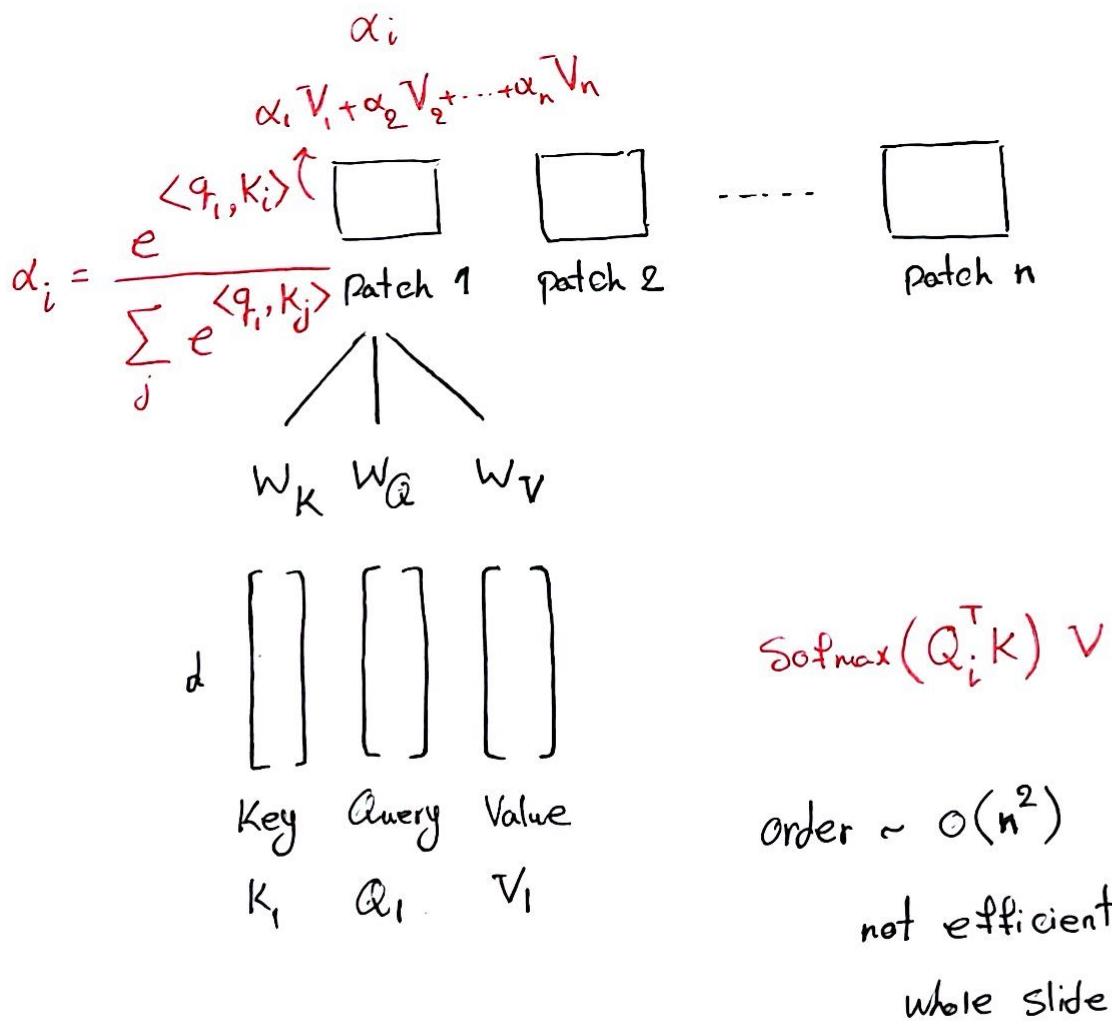
- aggregates embeddings then predicts label
- embedding-based → lack of interpretability
- instance-based → propagate error
- predicts each labels then aggregates them

Attention-based MIL

$$Z = \sum_{k=1}^K \alpha_k h_k \quad , \quad \alpha_k = \frac{e^{w^T \tanh(V h_k^T)}}{\sum_{j=1}^K e^{w^T \tanh(V h_j^T)}}$$

embedding feature map

Attention model



Dual-Stream MIL

like attention model but just calculate attention for main patch. q_m .

2 Streams { uses MIL max pooling for classification in first stream.
uses from q_f of main patch in classification to generate attention of other patches to q_m .

Graph Convolutional Network

$$h_i^{(l+1)} = \sigma \left(h_i^{(l)} w_0^{(l)} + \sum_{j \in N_i} \frac{1}{C_{i,j}} h_j^{(l)} w_1^{(l)} \right)$$

GNN / GCN

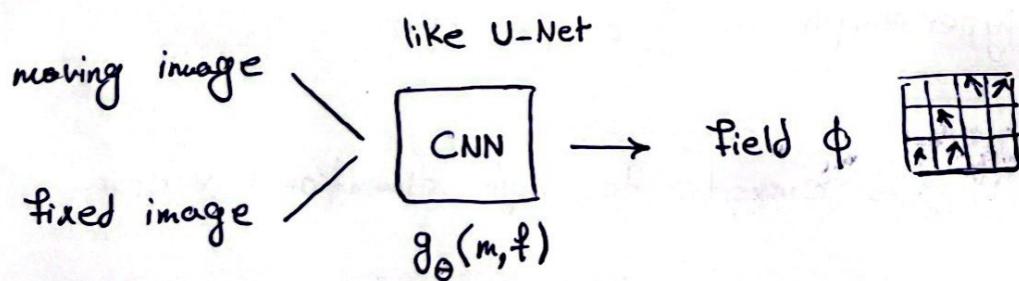
- Node classification
- Graph classification
- Link classification

Registration

- optimization problem for each 2 images / old approach
 - ⊗ $m \times n$ pixels ~ 9000 sec / each pair time consuming
 - ⊗ affine transformation (translation, Scaling, rotation)
it is not good when deformation exists.

$$\phi_{m,f} = \arg \min_{\phi} \underbrace{\| m \circ \phi - f \|}_{\text{image match}} + \underbrace{\lambda \text{Reg}(\phi)}_{\substack{\text{Smooth field} \\ (\text{regularization})}}$$

- it may give local min.
avoid from discontinuous and large field.



- unsupervised (Voxelmorph) $\sim L = \| m \circ \phi - f \| + \lambda \text{Reg}(\phi)$
 - loss
- fast for new image pair

$$p' = p + u(p) \sim \text{for being differentiable.}$$

Contour landmarks in image pair

- maximum dice score of these landmarks as loss or validation metric.

$$\text{Dice Score} = 2 \times \frac{\text{area of overlap}}{\text{area of Union}}$$

Synth Morph ~ uses from synthetic data

- for generalization of registration model.
 - Create imaginary data that are able to register.
 - OOD generalization.

Hyper Morph ~ $\underline{\lambda}$ as an input

- $\underline{\lambda}$ converts to high dimensional vector.
- Set an arbitrary $\underline{\lambda}$ in inference for custom dataset
- if variation of dataset is high you should set a small $\underline{\lambda}$ and vice versa.

Segmentation, Object detection, Instance Segmentation.

Semantic segmentation ~ U-net downsample & upsample

- up sampling approaches ~

unpooling {

- nearest neighbor $\begin{array}{|c|c|} \hline 2 & 3 \\ \hline 4 & 5 \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|c|} \hline 2 & 2 & 3 & 3 \\ \hline 2 & 2 & 3 & 3 \\ \hline \vdots & & & \\ \hline \end{array}$
- bed of nails $\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|c|} \hline 1 & 0 & 2 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \vdash & \vdash & \vdash & \vdash \\ \hline \end{array}$

Max unpooling ~ remember which element was max.

$$\begin{array}{|c|c|} \hline 2 & 3 \\ \hline 4 & 5 \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|c|} \hline 0 & 0 & (3) & 0 \\ \hline 0 & (2) & 0 & 0 \\ \hline \text{in previous position} & & & \\ \hline \end{array}$$

learnable - Convolution transpose ~ deconvolution

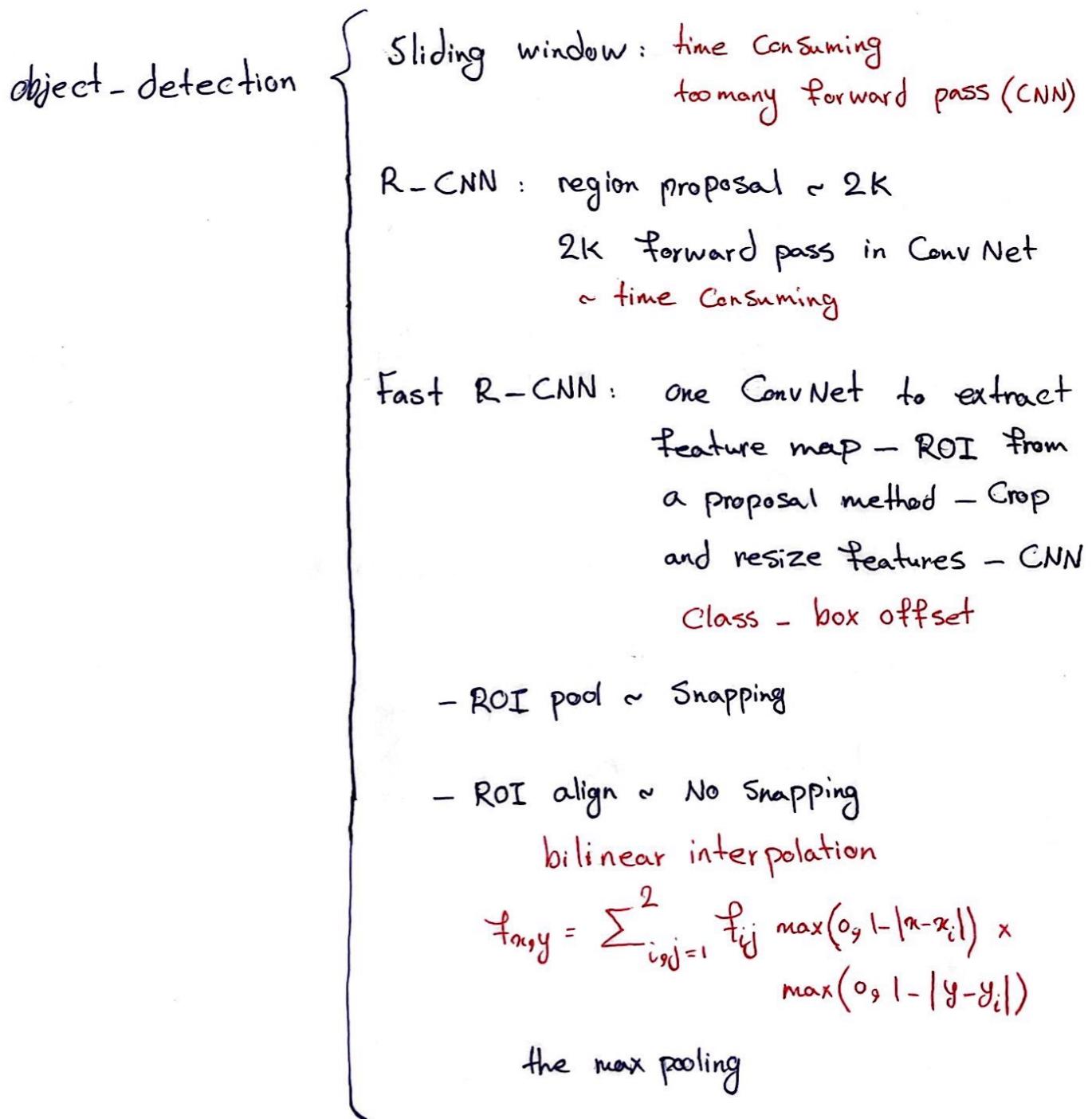
$$\left[\begin{array}{c|c|c|c} A_1 & A_2 & \dots & A_K \end{array} \right] \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{pmatrix} = \sum_{i=1}^K x_i A_i$$

* information loss is in latent layer (embeddings)

So we consider skip connection. \rightarrow in each step

why VNet needs small dataset to train?

- Fully Convolutional network has small number of weights
- 512×512 loss terms
- we can use a lot of transformation to augment data.
- High resolution convert to many patches.



the max pooling

Faster R-CNN

~ insert Region Proposal Network ~ $\frac{300}{\text{most Confident proposals}}$
for each pixel $\rightarrow K$ boxes
(in feature map)

ROI pooling \rightarrow Class, box offset

4 loss in training : $\begin{cases} \text{RPN} & \text{object/not object} \\ \text{RPN} & \text{regress box coordinates} \\ \text{Final} & \text{Classification} \\ \text{Final} & \text{box coordinates} \end{cases}$

Two Stage object detector $\begin{cases} 1. \text{ RPN} \\ 2. \text{ ROI pooling, class, box offset} \end{cases}$

Single-Stage object detector \rightarrow YOLO, SSD, RetinaNet

~ for each pixel $(5 \times B + C)$ \rightarrow (dx, dy, dh, dw, confidence)
aspect ratios \downarrow \rightarrow obj/not obj
 \downarrow classes

Mask R-CNN ~ like Faster R-CNN

~ with one more head for mask prediction.

- we should find also mask object in feature map
the ROI align it and enforce mask in mask prediction head.

Interpretability

→ assigns "attribution score" to each input pixel

- Attribution methods \sim Saliency map

$$\text{Saliency } (\alpha) := \nabla_{\alpha} f(\alpha) = \frac{\partial f(\alpha)}{\partial \alpha}$$

- Saliency maps are noisy!

Hypothesis 1: Saliency maps are truthful

noise is important to make decision

Hypothesis 2: gradients are discontinuous

DNN uses piece-wise linear functions
(ReLU activation, maxpooling, ...)

Sudden jumps in $f(\alpha)$ over small changes in input.

- Solution: smooth the gradient

$$\text{SmoothGrad}(\alpha) := \frac{1}{n} \int_1^n \frac{\partial f(\alpha^*)}{\partial \alpha^*}$$

$$\alpha^* = \alpha + N(0, \sigma^2)$$

gaussian noise

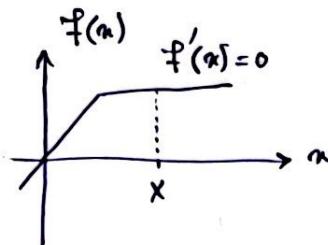
this solution smooth the jumps in $f(\alpha)$

debug
reliability
verify
discovery

right to explanation

Hypothesis 3: $f(u)$ saturates

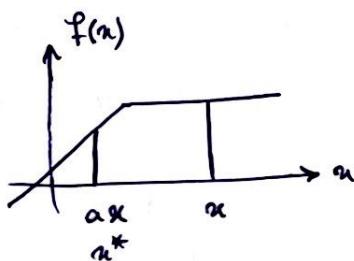
Small derivative locally. maximum possible activation.



but a feature may have strong effect globally.

- Solution: Interior gradient

scale back the x . $x^* = \alpha x \quad 0 < \alpha \leq 1$



$$\text{IntGrad}(x) := \frac{\partial f(x^*)}{\partial x^*}$$

Back prop - based

1. Deconvnet map feature pattern to input space

pass reconstructed signal through ReLUs.

removing noise by removing negative gradient

Backprop-based

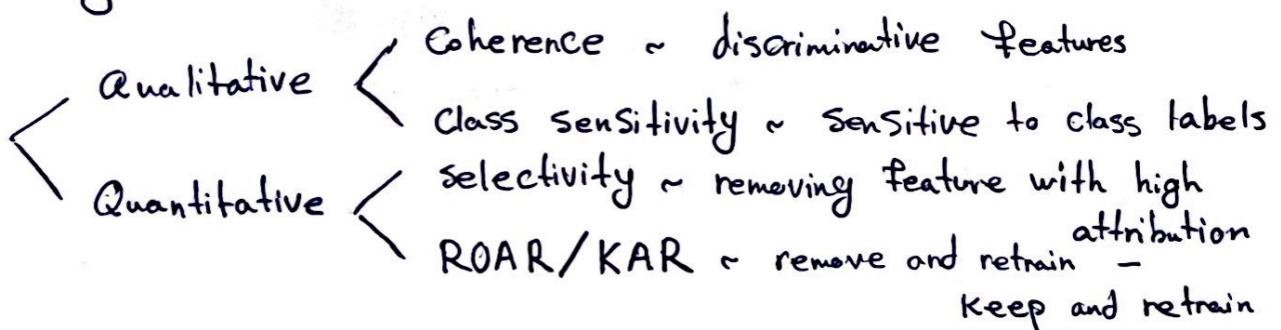
2. Guided backpropagation

Combine Deconvnet with back propagation

zero out the cells
which became zero
in forward pass.

removing negative gradient + consider forward activation

Evaluating Attribution methods



- selectivity ~ remove feature with high attribution
(salient feature) cause large decrease
in class probability.

if saliency features are very strong
we get a higher AOC (area over curve)

- ROAR / KAR ~ measure how the performance of classifier
Sensitivity may be not accurate changes as specific features are removed
based on attribution method.

ROAR: replace N% of pixels which are most important

KAR: replace N% of pixels which are least important

Federated learning

data cannot be centralized ~ motivation

- < too costly
- too sensitive ~ privacy

if we have centralized dataset, the use standard machine learning settings.

Each party learning on its own? No.

1. local dataset of each party may be too small.
 - overfitting
 - non-statistically significant result on medical data
2. local dataset may be biased.
 - not representative of target distribution.

Federate learning → collaboratively train a ML model while keeping the data decentralized

distributed - learning : data is centrally stored

- train faster
- data is distributed uniformly at random across workers.

FL → data is naturally distributed

- not independent and identically distributed
- imbalanced ~ total number of each dataset are not equal

Challenges in FL

- privacy constraints
- possibly limited reliability/availability of participants
- malicious parties & achieving robustness

FL

- Cross-device ~ massive number of parties $\leq 10^{10}$
 - ↳ has above challenges + small dataset per party
- Cross-silo ~ 2-100 parties
 - ↳ this is good

FL

- Server-orchestrated ~ server-client communications
 - global coordination, global aggregation
 - server is a single point of failure, bottleneck
- Fully decentralized ~ device-device communication
 - local aggregation
 - scales to a large number of devices
 - like graph

notation

set of K parties

party K ~ dataset D_K , n_K : number of data

$$D_1, D_2, \dots, D_K \quad n = \sum_K n_K$$

problem: $\min_{\theta} F(\theta; D) \xrightarrow{\text{parameter of model}}$

$$F(\theta; D) = \sum_{K=1}^K \frac{n_K}{n} F_K(\theta; D_K)$$

$$F_K(\theta, D_K) = \sum_{d \in D_K} f(\theta; d) \xrightarrow{\text{Cost function}}$$

} FedAvg
or
local SGD

Fed Avg

ρ : Client sampling rate

L : number of step in clientUpdate

$L=1, \rho=1 \rightarrow$ classic parallel SGD

updates are aggregated and model synchronized at each step.

* this algorithm is for server-orchestrated

FedAvg for fully decentralized setting

- undirected graph

$$G = (\{1, 2, \dots, K\}, E)$$

$\{k, l\} \in E \rightarrow k$ and l can exchange messages.

W matrix : weight 0 or 1 if $\{k, l\} \notin E$ then 0

$\Theta = \{\theta_1, \theta_2, \dots, \theta_K\}$ for each party

$$[W\Theta]_K = \sum_{l \in N_K} W_{K,l} \theta_l \rightarrow \begin{array}{l} \text{local aggregation} \\ \text{among neighbors of } \underline{K} \end{array}$$

↓ ↓
 for party K neighbors of K

Convergence rate depends on topology \sim more connected faster

two steps in each round

| | | |
|--|------|----------------------|
| 1. local update with $\theta_K^t - \frac{n_K}{B} \sum_{j \in B} \nabla f(\theta_K^t; d)$ | $\{$ | 2. local aggregation |
|--|------|----------------------|

Dealing with Non iid Data problem: Client Drift

Strong non-iid

- Scaffold

x^* is not average of x_1^* and x_2^*

- correction terms C_1, C_2, \dots, C_K for each party

update C as well as Θ in each round

- convergence rate is worse than parallel SGD

FedAvg slower than parallel SGD for strongly non iid data
~ large G

* learning from non-iid data is difficult because
each party wants the model to go in particular direction

Personalized model $\Theta_K \sim \text{Meta Model}$

- global model which easily adapts to each party

$$F(\Theta_1, \dots, \Theta_K; D) = \frac{1}{K} \sum_{k=1}^K F_k(\Theta_k; D_k) + \underbrace{\frac{\lambda}{2K} \sum_{k=1}^K \left\| \Theta_k - \frac{1}{K} \sum_{l=1}^K \Theta_l \right\|^2}_{\text{regularization}}$$

$$F(\Theta; D) = \frac{1}{K} \sum_{k=1}^K F_k(\Theta - \alpha \nabla F_k(\Theta); D_k)$$

ViT: Vision Transformers

- naive solution \rightarrow flatten image to height \times width vector
 16×16 image \rightarrow 256 vector
- transformer time complexity $\sim O(n^2)$ 256^2 too much!
- nice image representation
- Semi-supervised learning \sim linear classification
on iGPT representation

ViT is more general, they have less inductive bias.

CNN { locality
two dimensional neighborhood structure
translation equivariance

ViT { image patches 16×16 words
position embedding interpolation during fine-tuning + upscaling

→ { performs worse than CNN when there's not as much data
performs better than CNN when there's lots of data

inductive bias { ViT learns best representation for a task

- lack of inductive bias

CNNs are biased from beginning. they have inductive bias.

- locality \sim they generate feature using local pixels.

Global receptive field

- CNNs have a receptive field limited by kernel size.
- ViTs have global receptive field
- ViT shows distance between Q and K (two tokens)
mean attention distance is very high
when a token really want to attend
to another token.

ViT Cons

1. patches not very fine grained.
2. Some issues in task like segmentation, fine grained classification.
need exact boundary

You should scale up data and model size in Vision transformer, to get better performance.

* lack of inductive bias is bad when we don't have enough data.

Diffusion Model

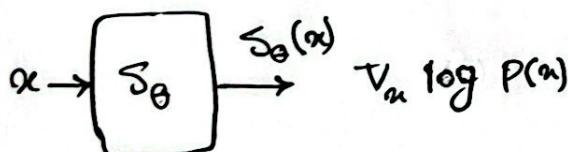
gradient of data distribution

$\nabla_{\theta} \log P(x)$: gradient ascent

Score function (matching)

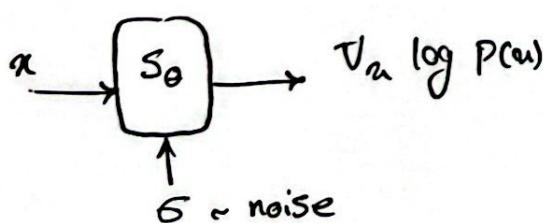
$$\mathbb{E}_{P(x)} \left[\text{trace} \left(V_m S_{\theta}(x) + \frac{1}{2} \| S_{\theta}(x) \|_2^2 \right) \right]$$

loss function for
training of S_{θ} network



$\mathbb{R}^{28 \times 28}$

non-sense gradient \sim Starting point is out of data manifold



Langevin dynamics (noise)

$$\tilde{x}_{t+1} \leftarrow \tilde{x}_t + \frac{\varepsilon}{2} S_{\theta}(\tilde{x}_t) + \sqrt{\varepsilon} z_t \quad z_t \sim N(0, I)$$

$$\sigma_1 > \sigma_2 > \dots > \sigma_k$$