

Ch1. Introduction



* DL (deep learning)



- f(x) complex, g(x) simple

① FCNet (fully connected)

② CNN (convolutional neural networks)

③ RNN (recurrent neural networks) : sequential data

- 딥러닝 1: CNN

Alex Net (5 convolution + 3 FC) \rightarrow VGGNet (deep CNN, 16+7 weight)

\rightarrow Google Net (22 layers) \rightarrow Deep Residual Net (1000+ layers)

- ML algorithm get simple

- 같은 dataset에 적용 가능. (performance↑)

- 딥러닝

① not sample-efficient. ② adversarial attack이 vulnerable

③ Suffer from hyperparameters.

* Problem set-up : $f(x) : \mathbf{x} \rightarrow \mathbf{y}$

Model

$$f(x) = w_1 \cdot \phi_1(x) + b_1$$

Error function

$$J = \sum_{n=1}^N (g(x_n) - y_n)^2$$

$$\frac{\partial J}{\partial w_1} = 2 \sum_{n=1}^N (g(x_n) - y_n) \phi_1(x_n) (-\phi_1'(x_n)) = 0$$

$$w_1 \left[\sum_{n=1}^N \phi_1^2(x_n) \right] = \left[\sum_{n=1}^N y_n \phi_1(x_n) \right]$$

* Least square method.

$$\sum_{n=1}^N (y_n - w_1 \phi_1(x_n))^2 = \sum_{n=1}^N e_n^2 = [e_1 \dots e_N] \begin{bmatrix} e_1 \\ \vdots \\ e_N \end{bmatrix} = e^T e = \|e\|^2$$

$$J(w) = \frac{1}{2N} \sum_{n=1}^N (y_n - w^T \phi(x_n))^2 = \frac{1}{2N} \|y - \phi^T w\|^2$$

$$\Rightarrow \text{极小值} : \text{이상 절 } \rightarrow \frac{1}{2N} \|y - \phi^T w\|^2 = \frac{1}{2N} (\mathbf{y} - \mathbf{w}^T \mathbf{y})^2 - \frac{1}{2N} \|\mathbf{w}\|^2$$

$$\frac{\partial}{\partial w} (\|y - \phi^T w\|^2) = \phi^T w - \mathbf{y} = 0$$

$$\phi^T w = \mathbf{y} \quad \therefore \mathbf{w}_L = (\phi \phi^T)^{-1} \mathbf{y} = \phi^T \mathbf{y}$$

* Linear model (likely hard to fit to randomness, y_n)

$$f(x_n) = w^T x_n + e_n \quad \text{BayesN} : E(y_n | x_n, \theta^*)$$

$$P(y|x) = N(w^T x + \theta^*, \sigma^2)$$

N samples (independent)

$$\prod_{n=1}^N P(y_n | w^T x_n) \stackrel{\text{log}}{\rightarrow} \sum_{n=1}^N \log P(y_n | w^T x_n)$$

* overfitting \Leftrightarrow weight vector w 이 \mathbb{B} 에 있는 것

- Ridge Regression

$$\min_w \frac{1}{2} \|y - \phi^T w\|^2 + \frac{\lambda}{2} \|w\|^2$$

$$\text{Lasso} : (\mathbb{E} \mathbb{E} + \lambda I)^{-1} \mathbb{E} y \quad (\lambda: \text{hyperparameter})$$



* logistic model.

$$f(x) = P(y=1|x) = \frac{1}{1 + e^{-w^T x}}$$

$$p(y|x) = \text{Bernoulli}$$

$$p(y|x) = P(y=1|x)^y P(y=0|x)^{1-y}$$

$$-\text{objective} : \sum_{n=1}^N p(y_n|x_n)$$

- error function $J(w)$

$$\text{① Gradient Descent} : w^{(k+1)} \leftarrow w^{(k)} - \alpha \nabla J(w^{(k)})$$

$$\text{② Newton's method} : w^{(k+1)} \leftarrow w^{(k)} - \alpha H^{-1} \nabla J(w^{(k)})$$

$$\text{Hessian} = \nabla^2 J(w)$$

* softmax regression.

- multi class (ex. 3)

$$\text{compute } p(y=1|x) \quad w_1^T x \\ p(y=2|x) \quad w_2^T x \\ p(y=3|x) \quad w_3^T x \quad \Rightarrow \quad P(y=k|x) = e^{w_k^T x}$$

$$- P(y_n = k|x_n) = \text{softmax}(Q_n) = \frac{e^{w_k^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$\text{ex. } (\text{Model}) \rightarrow \begin{cases} j &= 0, 1 \\ P &= \begin{bmatrix} 0.1 & 0.2 \\ 0.1 & 0.1 \end{bmatrix} \end{cases} \quad \begin{cases} j &= 0 \\ &1 \\ &2 \\ &0 \\ &1 \\ &0 \end{cases} \quad \begin{cases} \text{Tiger} \\ \text{lion} \\ \text{cat} \\ \text{dog} \\ \text{bird} \\ \text{horse} \end{cases} \quad \begin{cases} \text{time: 0: False} \\ \text{1: True} \end{cases}$$

$$\text{cross entropy } H(p, t) = -\sum_j p_j \log t_j$$

$$CE = -1 \cdot \log 0.1 - 0 \cdot \log 0.2 - 0 \cdot \log 0.1 = -\log 0.1 = 0.31$$



Ex. Suppose

$$D = \{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}, y_i = w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$y_1 = w_1 \cdot 1 + w_2 \cdot 1 + w_3 \cdot 0$$

$$y_2 = w_1 \cdot 1 + w_2 \cdot 0 + w_3 \cdot 1$$

$$y_3 = w_1 \cdot 0 + w_2 \cdot 1 + w_3 \cdot 1$$

determine w_1, w_2, w_3 (L_2 loss + L_1 regularization).

① # equations, ② # unknowns

① > ② : over-determined, under-complete \Rightarrow least-squared regression

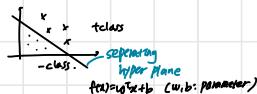
② > ① : Complete \Rightarrow Ridge

③ < ② : over complete \Rightarrow LASSO

ch3. Perception

- perception: single layer neural network (hidden layer x)
- multilayer perception: multi layer extension of perception.
- = Fully connected Network.

* linear classification



XOR function

$$J(w) = \sum_i w_i x_i j_i$$

$$w = \arg \min_w \left[-\frac{1}{n} \sum_{i=1}^n w_i x_i j_i \right] \text{ apply gradient descent}$$

$$\text{Gradient: } \nabla J(w) = x_i j_i$$

* Algorithm outline

- ① train sample x_n
- ② misclassified x_n check
→ if nothing
→ if yes: update $w \leftarrow w + w_n = w + x_n j_n$

* AND operation



* XOR



* McCulloch-Pitts Model



$$z_1 \rightarrow O \rightarrow o_1 \quad z_2 \rightarrow O \rightarrow o_2 \quad \text{activation function}$$

$$z_3 \rightarrow O \rightarrow o_3 \quad z_4 \rightarrow O \rightarrow o_4 \quad \text{function}$$

$$h = p(z_1) + q(z_2)$$

$$h = q(z_3) + r(z_4)$$

$$j = V_1 h_1 + V_2 h_2 = V_1 p(z_1) + V_2 q(z_2)$$

$$= V_1 p(z_1) + V_2 q(z_2) \quad \text{adaptive bars}$$

$$\text{linear model: } j = V_1 p(z_1) + V_2 q(z_2) \rightarrow \text{basis fixed.}$$

* activation function

$$\text{Sigmoid: } y = \frac{1}{1+e^{-x}} \quad \text{② tanh, } y = \tanh(x) \quad \text{pre-activation}$$

$$\text{③ ReLU: } y = \max(0, x) \quad \text{④ Leaky ReLU: } y = \max(0, \alpha x, x)$$

→ 경계값에 미분 불속 → subgradient

$$\frac{\partial L}{\partial w} = \frac{1}{n} \sum_{i=1}^n \delta_i x_i$$

* RelU의 강점

① cheap

② converge faster

③ sparsely activated



* softplus

$$y(x) = \ln(1+e^x)$$

* MLP structure



* back propagation

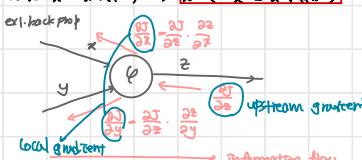
$$\begin{aligned} & \text{- gradient descent} \\ & \min_w J(w) \end{aligned}$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - f_w(x_i))^2$$

$$\frac{\partial J}{\partial w} = \frac{1}{n} \sum_{i=1}^n (2f_w(x_i) - 2x_i) = 2f_w(x_i) - 2x_i = 0$$

$$\therefore x_i = x_i - \frac{1}{n} \sum_{i=1}^n (2f_w(x_i) - 2x_i)$$

gradient \equiv weight update



Ex.2. $x \rightarrow w_1 \rightarrow w_2 \rightarrow y$

$\frac{\partial J}{\partial w_2}$: directly calculated

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial w_2} \Rightarrow \text{Chain rule 적용.} \rightarrow \text{chain rule 적용.}$$

$$\begin{aligned} \frac{\partial J}{\partial w_2} &= \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial w_2} \\ &= \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial w_1} \cdot \frac{\partial w_1}{\partial w_2} \\ &= \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial w_1} \cdot \frac{\partial w_1}{\partial w_2} \end{aligned}$$

Ex. Calculate Gradient by back prop.

$$\frac{\partial J(w)}{\partial w} = \frac{1}{1+e^{C(w_1 + w_2 x_1 + b)}}$$

$$\begin{aligned} w_1 &= 0.1 \\ w_2 &= 0.2 \\ x_1 &= 0.4 \\ x_2 &= 0.6 \\ b &= 0.1 \end{aligned}$$

ch4. Optimization

$$\begin{aligned} & \text{- error function} \\ & J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2 \quad | \theta: \text{model parameters} \quad \text{Square loss: } J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2 \end{aligned}$$

$$\Rightarrow \text{Want: } \theta = \arg \min J(\theta), \text{ error } \Rightarrow \text{loss function}$$

* Gradient Descent: iteration 반복하면서, 그때마다 gradient를 계산해 step size를

$$\theta_{\text{new}} \leftarrow \theta_{\text{old}} - \left[\frac{\partial J(\theta)}{\partial \theta} \right] \quad \theta_{\text{old}} = \theta_{\text{old}} + dt \frac{\partial J(\theta)}{\partial \theta} \quad \text{current gradient}$$

Gradient = steepest descent

* Batch & mini-Batch



$$\text{full batch: } J(\theta; X; Y) = \frac{1}{m} \sum_{i=1}^m (y_i - f_\theta(x_i))^2$$

→ training dataset 전체 이용. :: memory 문제를 발생

$$\text{mini-batch: } \theta_{\text{new}} \leftarrow \theta_{\text{old}} - \frac{1}{B} \sum_{i=1}^B (y_i - f_\theta(x_i))^2$$

→ training batch에 대해서만 gradient 계산
if batch size가 1이라면 sample by sample 사용.

* Saddle point

! optimization 문제로 자주 만날 때 min or max 둘 다 (3 dimension 이상일 때)

* Stochastic Gradient Descent = SGD

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2 \approx E[y - f_\theta(x)]^2$$

$$-\theta \leftarrow \theta - d\theta \cdot (y - f_\theta(x))^2 \quad \text{batch size이 } B.$$

$$d\theta = \frac{1}{B} \sum_{i=1}^B (y_i - f_\theta(x_i))^2$$

$$d\theta = \frac{1}{B} \sum_{i=1}^B (y_i - f_\theta(x_i))^2$$

$$\Rightarrow J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2 + \frac{1}{B} \sum_{i=1}^B (y_i - f_\theta(x_i))^2$$

mini-batch Gradient descent

$$= SGD$$

→ problem: error (정답) noisy!!

step size가 정답 찾기 오류.

Ch9. CNN (Convolutional Neural Network)

* FC-nets

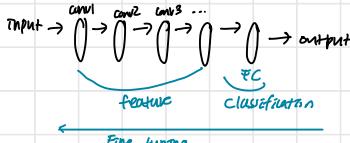
$$\text{Input} \rightarrow \begin{matrix} h_1 \\ h_2 \\ h_3 \\ h_4 \end{matrix} \rightarrow \text{output} = h_1 = f(a_1, a_2, a_3) = (w_{11}a_1 + w_{12}a_2 + w_{13}a_3)$$

* CNN

$$\text{Input} \xrightarrow{\text{hidden}} \begin{matrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \end{matrix} \rightarrow \text{output} = h_1 = f(w_1x_1 + w_2x_2 + w_3x_3) \\ h_2 = f(w_4x_1 + w_5x_2 + w_6x_3 + w_7x_4) \\ h_3 = f(w_8x_1 + w_9x_2 + w_{10}x_3 + w_{11}x_4) \\ h_4 = f(w_{12}x_1 + w_{13}x_2 + w_{14}x_3 + w_{15}x_4) \\ h_5 = f(w_{16}x_1 + w_{17}x_2 + w_{18}x_3 + w_{19}x_4)$$

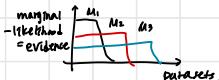
Convolution \Rightarrow cross-correlation between $\begin{bmatrix} x_1 & x_2 & x_3 \\ x_4 & x_5 & x_6 \end{bmatrix}$ and $\begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \end{bmatrix}$

- parameter sharing



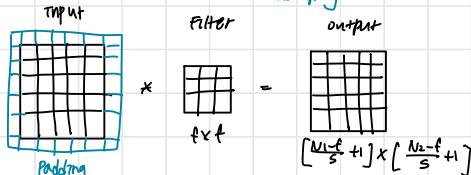
* Bayesian Model comparison

$$- \text{Model selection} : P(D|M_i) = \int p(D|w_i, M_i) p(w_i|w_i) dw_i$$

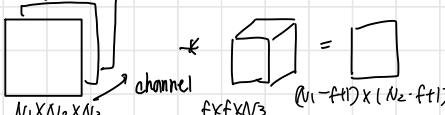


* Convolution

- Input layer \rightarrow feature maps \rightarrow pooling \rightarrow convolution \rightarrow pooling \rightarrow features
 \hookrightarrow downsampling

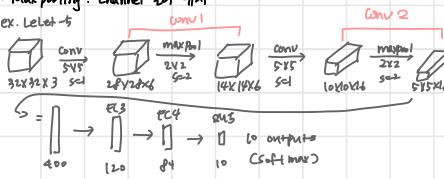


- Convolution over volume



- Channel of 32는 32개의 1x1 convolution 결과.

- Max pooling : channel 32로 풀기



* UpSampling

\hookrightarrow Unpooling & Transposed Convolution

① Nearest Neighbor

\rightarrow 2x2 in 2x2

Ex. $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ upsampled

Input 2x2

a	b
c	d

Output 4x4

② bilinear

\rightarrow max unpooling

\rightarrow select max that neighbor 가 있는 4부분

- Transposed Convolution

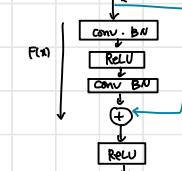
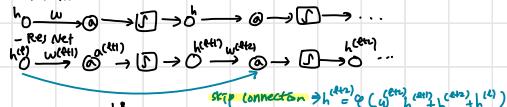


$$\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} w_1a + w_3c & w_1b + w_3d \\ w_2a + w_4c & w_2b + w_4d \end{bmatrix}$$

- filtered 가중치는 그대로 frontole는 그대로 \Rightarrow $\frac{1}{2}$

* ResNet

- Plain net



$$- \text{gradient } \frac{\partial J}{\partial x} = \frac{\partial J}{\partial F} \frac{\partial F}{\partial x} + \frac{\partial J}{\partial x}$$

* Inception Nets

: loss는 통상으로 나누어 계산하기 때문 total loss를 찾는다

- V1, V3, V4 ...

* object detection

= class prediction + Localization

- R-CNN : object가 있는지 없는지 판별을 \Rightarrow propose \hookrightarrow detect

\rightarrow 2阶段 framework \Rightarrow 디비전 및 추적

Ex: AP ↑ (object detection accuracy)

- fast R-CNN : Faster CNN만 사용 \Rightarrow speed↑ 비용↑ but overhead↑

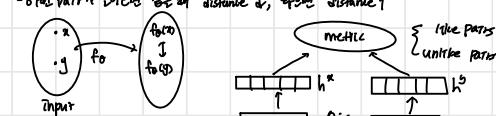
- faster R-CNN : Region proposal \Rightarrow overhead 줄임, 비용↑

- YOLO : grid 344x344에 9000개의 bounding box를 통해 한 번에 처리

* CNN with Share Architecture

\Rightarrow Detection & Verification : 같은 모델에서 같은 두 모양을 판별

- total parallel branches \Rightarrow distance d_1 , d_2 , ..., distance d_n



* CNN for timeseries classification

Ex: 웹사이트 방문하는 패턴



- pre-trained CNNs

① VGG

② Inception Net (1x1x3, 5x5x3, maxpool) \rightarrow 3層 deep CNN

③ ResNet (skip connection)

$$h^{(L)} \xrightarrow{\text{Conv}} \xrightarrow{\text{BN}} \xrightarrow{\text{ReLU}} \xrightarrow{\text{Conv}} h^{(L+1)} \xrightarrow{\text{ReLU}}$$

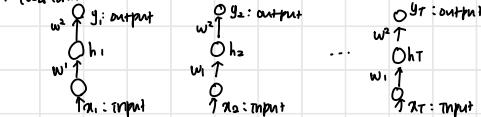
- Image classification

- semantic segmentation (upsampling) \rightarrow Verification (Share Architecture)

- object detection (YOLO)

Ch8. RNN (Recurrent Neural Network)

*Feed forward network



$$[x_1] \rightarrow \text{model} \rightarrow [y_1]$$

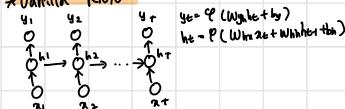
if input sequence (ex. text, video, music, ...)

$$[x_1, x_2, \dots, x_T] \rightarrow \text{model} \rightarrow [y_1, y_2, \dots, y_T]$$

→ 각각을 independent 히든 상태로 처리 (dependency 고려 필요)

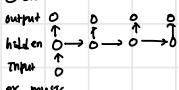
∴ hidden transition 필요

*Vanilla RNN

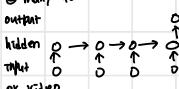


→ Variable x_t 는 handling ht-1과 xt이다

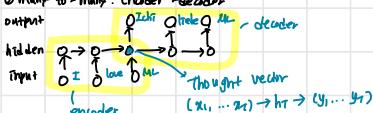
① One-to-many



② many-to-one



③ many-to-many: encoder-decoder



→ seq-to-seq learning

- backprop through time

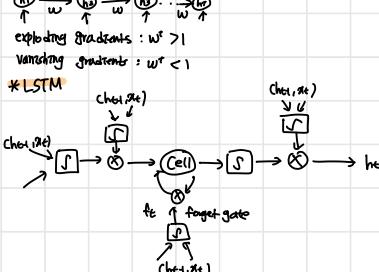
$$\text{loss} = \text{loss}_1 + \text{loss}_2 + \text{loss}_3 + \text{loss}_4 + \text{loss}_5 \rightarrow \text{total loss.}$$



exploding gradients: $w^t > 1$

vanning gradients: $w^t < 1$

→ LSTM



- forget gate: 0 or 1. If forget gate = 0 : cell의 초기 상태 사용

if forget gate = 1 : forget gate를 적용해 초기화

→ input-output, forget gate는 sigmoid 활성화 함수로 처리되어 0과 1 사이의 값을 가짐

- gradient flow



*GRU

· cell state update gate z_t 및 reset gate r_t 로 구현되는 구조

[If $z_t = 1$: $h_t = g_t$ (ht는 초기)]

[If $z_t = 0$: $h_t = h_{t-1}$ (ht는 초기)]

$$h_t = (1-z_t) \odot h_{t-1} + z_t g_t$$

- natural language에서 GRU보다 사용. ordered sequential data는 LSTM보다 사용

*Generative RNN

$$(y_1, p(y_1|x_1)) \rightarrow (y_2, p(y_2|x_1, x_2)) \rightarrow (y_3, p(y_3|x_1, x_2, x_3))$$

→ 각각은 independent 히든 상태로 처리됨 (dependency 고려 필요)

∴ hidden transition 필요

K Teacher-forcing

$$g_t \rightarrow h_t \rightarrow y_t$$

CTF

: Set y_t 의 true value

→ training은 각 차례별 test N 문제

→ scheduled sampling, 폴

→ scheduled sampling.

→ training 때 true value y_t 를 input으로 넣음 (TF) → 예상 결과

→ 예상 결과와 원래의 y_t를 비교해 평가됨. (teacher-learning)

Ch9. Attention & Transformer model

*Attention in Encoder - Decoder

: encoder - decoder에서 output을 제작하는 데 사용되는 encoder의 input을 활용해 이를 활용해 출력을 생성하는 방식

- 초기값은 context vector = $\sum_i^t \text{dec}_i \cdot \text{enc}_i$ → context sum, Zeros

- dec: output y_t 의 input을 얻기 위해 이 attention mechanism을 찾음

$$\Rightarrow \text{dec}_t = \frac{\sum_{i=1}^{t-1} \text{enc}_i}{\sum_{i=1}^{t-1} \text{enc}_i}$$

- Alignment model: input의 어느 부분이 더 잘 대응하는지 찾음

$$\text{Seq} \rightarrow \text{FF} \rightarrow \text{enc}_t$$

$h_t \rightarrow \text{enc}_t$ (se: hidden node in decoder, enc: encoder의 hidden node, seq: sentence)

→ input sequence와 output sequence의 alignment는 1이거나 0

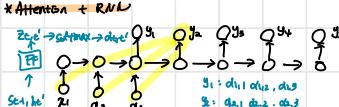
*Visual Attention

: input image를读懂할 때 output을 찾음.

- encoder: CNN 사용, decoder: RNN 사용

→ attention을 통해 image와 문장과 같은 패턴을 찾음

*Attention + RNN



*Transformer models

O O O ... O vector token

$x_1, x_2, x_3 \rightarrow \text{RNN}$: sequential temporal dependence. (RNN 형태)

CNN: 영상 layer의 패턴 (영상 맵 → 패턴 맵)

- sequence를 활용한 RNN, CNN이 아닌 attention만 사용

*Vanilla transformer

: encoder + decoder = attention + positional encoding + feedforward net

- dot product attention

$$\text{query } q_i \rightarrow \sum_j w_{ij} v_j, w_{ij} = \frac{e^{q_i^T k_j}}{\sum_l e^{q_i^T k_l}}$$

Input seq: (x_1, \dots, x_n) Output seq: (y_1, \dots, y_n)
 $q_i = \sum_j w_{ij} x_j$ ex. $y_1 = \sum_i w_{1i} x_i + \sum_i w_{2i} x_i + \sum_i w_{3i} x_i$, $w_{1i} = \frac{e^{q_1^T k_i}}{\sum_l e^{q_1^T k_l}}$

→ query: 단어에 대한 가치지

- key: 단어에 대한 쿼리에 대한 가치지

- value: 단어에 대한 가치지

- scaled dot-product attention

$N = \text{seq. length}$

$$g_{ij} = \frac{1}{\sqrt{d_k}} \left(\frac{e^{\frac{q_i^T k_j}{\sqrt{d_k}}}}{\text{softmax}(\frac{q_i^T k_j}{\sqrt{d_k}})} \right) v_n$$

- multi-head attention

head i : $\text{Attention}(Q_i w_i^Q, K_i w_i^K, V_i w_i^V)$

- encoder: self-attention layers.

- decoder: self-attention layers, input \rightarrow masked

- Positional encoding

① 1st time step initial unique encoding

② 2nd time step initial learned consistency

③ 3rd sequence initial randomization

④ deterministic

ex. 0 : 0 0 0 0 0

1 : 0 0 0 1 0

2 : 0 0 1 0 0

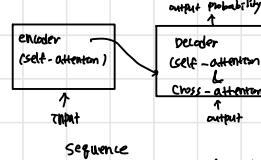
3 : 0 0 1 1 0

4 : 0 1 0 0 0

5 : 0 1 0 1 0

→ $\text{pos. embed.} + \text{vector} = \text{Positional vector} + \text{embedding vector} \Rightarrow \text{structural}$

→ Position information \Rightarrow seq.



→ need to consider "context"

할 수 있는 것은 y_1, y_2, \dots, y_t 를 통해 y_i 를 찾는 것.

$$x_1, x_2, \dots, x_t \rightarrow \text{self-attention} \rightarrow y_1, y_2, \dots, y_t$$

→ 각각의 단어가 다른 문장에 등장하는 위치에서 찾을 수 있는 단어들이다.

→ 예상 단어를 찾기 위해 단어의 \rightarrow encoding \rightarrow 찾기 가능

* Transformer vs Reformer

- Transformer: long sequence X , dot product attention needs complexity

- Reformer: LSH attention, reversible layers, chunking

↳ dot product attention \rightarrow complexity ↓

↳ memory space ↓

* Hamming space (Locality Sensitive Hashing)

- Hamming space: $[0, 1]^n, \dots, [1, 0]^n$

→ high dimension data $\rightarrow 0.1$ 차이를 두 dimension 씩 0.125 \rightarrow 3D

- probability collision: 2nd-level vector \rightarrow LSH \rightarrow ↑

$$\begin{array}{c} \text{① } \begin{array}{c} \oplus \\ \text{Hamming distance} \end{array} \\ \oplus \quad \oplus \quad \oplus \quad \oplus \quad \oplus \\ \oplus \quad \oplus \quad \oplus \quad \oplus \quad \oplus \end{array} \rightarrow \text{Hamming distance} \rightarrow \text{collision rate} \rightarrow \text{probability of collision}$$

→ probability collision: 2nd-level vector \rightarrow LSH \rightarrow ↑

$$x_1 \rightarrow y_1 \rightarrow z_1 \rightarrow z_2 \rightarrow z_3 \quad z_1 = y_1 + \epsilon_1(y_1)$$

$$z_1 \rightarrow z_2 \rightarrow z_3 \quad z_2 = z_1 + \epsilon_2(z_1) \quad z_3 = z_2 + \epsilon_3(z_2)$$

$$y_1 = \text{ReLU}(L(x_1)) \quad z_1 = \text{ReLU}(L(z_1))$$

* Reversible network

- Random projection은 정부한 sequence \rightarrow

* Permutation invariant

- Input vector

$$[] \rightarrow \text{model} \rightarrow \text{output vector} \rightarrow []$$

- Input sequence

$$[], [], \dots, [] \rightarrow \text{model} \rightarrow \text{output sequence} \rightarrow [], [], \dots, []$$

- Input set

$$\{ [], [], \dots, [] \} \rightarrow \text{model} \rightarrow []$$

↳ set의 원소의 순서를 바꾸어도 같은 set

↳ set의 원소의 순서를 바꾸어도 같은 set

↳ permutation-invariant 필요

* Deep sets

- 기계학적 permutation invariant

$f(x) \approx P(\text{set}(x), \theta)$

* Permutation Invariance and Equivariance

- permutation invariant: set x 의 모든 원소를 바꾸어도 같은 response를 갖다

- permutation equivariant: mapping permutation하는 layer outputs permutation

→ 순서도 동일하게 바꾸어 준다

$$\text{ex. } N=3, \theta = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} + \sigma \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0.1 & 0.2 & 0.3 \\ 0.2 & 0.3 & 0.1 \\ 0.3 & 0.1 & 0.2 \end{pmatrix}$$

* Set transformer model



- permutation-equivariant Attention modules (SAB & ISAB)

↳ query set \rightarrow key \rightarrow self-attention (auto encoder style)

- Pooling by Multi-head attention

- encoder: stack of SAB & ISAB

- set transformer's permutation invariant

↳ SAB은 encoder의 permutation-equivariant인 layer로 인해

→ multi-head attention은 pooling

→ aggregation, set \rightarrow global feature vector를 갖는 (permutation-invariant)

→ Decoder \rightarrow target's mapping

* Set transformer examples

① Meta Anomaly Detection

② Point cloud classification

↳ Shape datasets 학습 및 classification

Ch10. Deep Generative models

* Discriminative models



* Density Estimation

P(model($x; \theta$)) \approx P(data(x))

↳ likelihood \rightarrow empirical distribution

- Prescribed models

① Fit $P_\theta(x)$ (model distribution) to $P_\text{data}(x)$ (empirical distribution)

→ parameterized distribution likelihood (특정 분포에 대한 확률)

② assign probability to every x in data distribution

ex. Variational autoencoder (VAE)

- Implicit models

① learn generative network, $z \rightarrow g(z)$

→ likelihood \rightarrow (continuous setting)

ex. GAN (generative adversarial networks)

- choose KL-divergence

$$\theta = \arg\min_{\theta} D_\text{KL}(P_\theta(x) \parallel P_\text{data}(x))$$

$$D_\text{KL}(P_\theta(x) \parallel P_\text{data}(x)) = \int P_\text{data}(x) \log \frac{P_\theta(x)}{P_\text{data}(x)} dx = \int P_\text{data}(x) \log \frac{1}{P_\theta(x)} dx$$

$$\Rightarrow \arg\max_{\theta} \int P_\text{data}(x) \log \frac{P_\theta(x)}{P_\text{data}(x)} dx \Rightarrow \arg\max_{\theta} \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log \frac{P_\text{data}(x)}{P_\theta(x)} dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log \frac{P_\text{data}(x)}{P_\theta(x)} dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log \frac{P_\text{data}(x)}{P_\theta(x)} dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

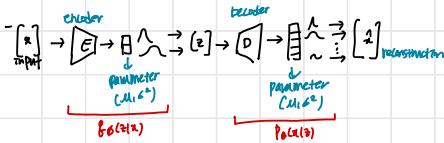
$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

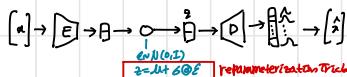
$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\text{data}(x) dx$$

$$\Rightarrow 0 + 0 + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx$$

$$\Rightarrow \int P_\text{data}(x) \log P_\theta(x) dx - \int P_\text{data}(x) \log P_\theta(x) dx + \int P_\text{data}(x) \log P_\text{data}(x) dx - \int P_\text{data}(x$$



* Train VAE with Reparameterization Trick



- Consider a random variable $z = \mathcal{Z} \sim N(\mu, \sigma^2)$
- $\mathbb{E}[z] = \mu$, $\text{Var}(z) = \sigma^2$
- $z = \mu + \sigma \mathcal{E}$, $\text{Var}(z) = \sigma^2$

* Training VAE - maximizing ELBO (Evidence lower bound)

- Variational lower-bound

: log-likelihood

$$\log P_\theta(x) = \log \int p_\theta(x|z) dz = \log \left\{ \int p_\theta(z|x) \frac{p_\theta(x|z)}{p_\theta(z|x)} dz \right\} \geq \mathbb{E}_{p_\theta(z|x)} \left[\log \frac{p_\theta(x|z)}{p_\theta(z|x)} \right] = \mathbb{E}_{p_\theta(z|x)} \left[\log p_\theta(x|z) \right] - \mathbb{E}_{p_\theta(z|x)} \left[\log \frac{p_\theta(z|x)}{p_\theta(z)} \right]$$

↳ Jensen's inequality
↳ Reconstructor ↳ Penalty

* Noisy Gradients

$$\begin{aligned} - \nabla_{\theta} \mathbb{E}_\theta [\log p_\theta(x|z)] &= \nabla_{\theta} \int p_\theta(z|x) \log p_\theta(x|z) dz \\ &= \int \nabla_{\theta} p_\theta(z|x) \cdot 0 \cdot p_\theta(z|x) \frac{p_\theta(x|z)}{p_\theta(z|x)} dz = \int \nabla_{\theta} p_\theta(z|x) \log p_\theta(x|z) p_\theta(z|x) dz \\ &= \mathbb{E}_\theta [\nabla_{\theta} p_\theta(z|x) \nabla_{\theta} \log p_\theta(x|z)] \end{aligned}$$

* Regularized Autoencoder (RAE)

- No noise injection

- deterministic autoencoder + explicit regularization for the decoder

- loss for RAE : $J_{RAE} = J_{RE} + \lambda J_{REG}$
 \downarrow \downarrow \downarrow
 ↳ Explicit regularizer for the decoder

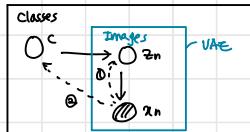
- Example of JREG

Distributional regularizer : $J_{REG} = \frac{1}{2} \|B\|^2$ (weight decay on decoder B)

② Gradient penalty : $J_{REG} = \|\nabla D_\theta(E_\theta(z))\|^2$

③ Spectral normalization

④ Neural statistician



- ① Standard inference network (VAE-like)
- ② State network (permutation-invariant model)

Ch11. Generative Adversarial Networks (GANs)

* Adversarial Training



- GAN이 이상의 각각의 충돌하는 criterions을 갖는 모델들이 equilibrium을 이루며 GAN을 학습함.

* Training GANs

- Generator : $\min_{G} [\mathbb{E}_{\text{Empirical}} [\log (1 - D(G(z; \theta)))]]$

- Discriminator : Generator의 출력을 갖춘 %, $\alpha \rightarrow (D) \rightarrow \text{P}(\text{real}) = 1 - \text{P}(\text{fake})$
 ↳ 허위 생성물

$$\rightarrow \max_{D} [\mathbb{E}_{\text{Empirical}} [\log D(x; \theta)] + \mathbb{E}_{\text{Empirical}} [\log (1 - D(G(z; \theta)))]]$$

* Generating Images by GANs

- DC GAN : using 100 dimensional random noise vector
- Progressive Growing of GANs : speed up + less noise
- GANs for style image super-resolution : upscaling을 사용해 더 선명함

* GANs vs. AE

① Non Convergence

② Model collapsing

: generator가 discriminator를 속여버려 반복됨
 → entangled mode : discriminator를 잘 쓸 때 model은 사용할 수 없음 (48x)

③ diminished gradient

: discriminator는 generated gradient를 처리함

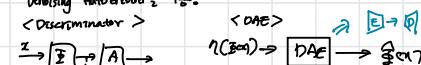
* Decreasing Autoencoder



noise contaminated channel $\xrightarrow{\text{error}} \text{clear: } x - \hat{x}$

* GANs trained with Denoising Feature matching

- Denoising Autoencoder는 48x:



*

	VAE	GAN
training	reg	(iterative) training
image	little blurry	sharp looking

* Semi-Supervised learning

- labeled example 1/2, unlabeled example 2/2에 의해

- loss : Lsupervised + Lunsupervised

↳ GAN + GAE

* Conditional GAN

- generate a fake sample with a condition y

→ generator는 condition에 의해

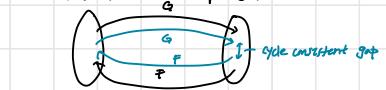


- ex: unpaired Images to Image translation

(A가 B와 같은 특징을 가지고 있음)

→ Different pairs of images have different domains

Photo X Photo Y



: paired GANs

* Adversarially learned Interference

- Inference net $\xrightarrow{(x, z)}$ discriminator $\xrightarrow{(x, z)}$ $\xrightarrow{\text{criterion}}$

- ↳ GAN은 Decoder는 GAN이

Ch12. Bayesian Optimization

* Black box optimization

: black box function -> not differentiable or f(x) & g(x) query cost
x: cost & f(x) (or g(x) may)

* Regret

- instantaneous regret : better chosen than

- cumulative regret

- simple regret

- cumulative regret: linearly related to simple regret

* Hyper parameter

- regularize hyper parameter (useless may) \rightarrow ~~useless may~~

ex. learning rate, batch size, kernel size, # of layer...

- model parameter

: Data를 통한 모델의 weight %

ex. weight

- Hyper parameter은 빠른 %

① grid search

② random search

③ bayesian optimization

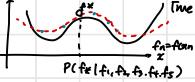
* Auto ML

- Feature processing

- Model / Algorithm selection

- Hyperparameter tuning

* Gaussian Process Regression



- 장점

① good uncertainty estimates

② choice of reasonable kernel can approximate binned data \rightarrow ~~etc.~~

- 쓰임

① Not scalable

② limited by kernels

- kernel \rightarrow ~~etc.~~

① RBF kernel (squared exponential kernel)

② Matern kernel

GPR: decision tree, Neural network \neq GP

- $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$

$$\text{Posterior mean } \hat{\mathbf{f}}(x) = \mathbf{C}(x, \mathbf{x}_n) [\mathbf{C}(\mathbf{x}_n, \mathbf{x}_n) + \sigma^2 \mathbf{I}]^{-1} \mathbf{y}_n$$

$\downarrow \frac{\partial}{\partial \mathbf{x}_n} \mathbf{C}(\mathbf{x}_n, \mathbf{x}_n)$

$\downarrow \frac{\partial}{\partial \mathbf{x}_n} \mathbf{y}_n$

* Variance of Prediction (Uncertainty U)

* Random forest

- 빠르고 정확한 %

- Poor extrapolation

* Neural Network

- Steepest Descent: ~~etc.~~

* Categorical / Integer - Valued Variable

- Categorical: One-hot encoding

- Integer - valued: Rounding after optimizing the acquisition function

* Acquisition function

: optimization 목적의 주제는 대체로 최적화 문제를 풀 때나 그 외에 다른

input 변수를 최적화하는 %

① Probability of Improvement (PI)

② Expected Improvement (EI)

③ Global Confidence Bound (GP-UCB)

④ EI : $U(x, y; \theta) = L(y; \theta)$

$$\rightarrow y \in \text{distr. } \theta, \text{ where } \theta$$

$P(x) = P(f(x) \geq f(x_0))$

- ~~etc.~~ : slack variable θ 가 %

\rightarrow too big θ : slow tuning, too small: too narrow

* EI

$$EI(x) = \max(y - \mu, 0) \sqrt{L(y; \theta)} dy$$

- PI와 비슷한데 ~~etc.~~

$$\therefore EI(x) = \begin{cases} (L(x) - \mu(x)) \sqrt{L(x; \theta)} & \text{if } L(x) > 0 \\ 0 & \text{if } L(x) = 0 \end{cases}$$

$$\theta = \begin{cases} \frac{\max(L(x)) - \mu}{\sigma} & \text{if } L(x) > 0 \\ 0 & \text{if } L(x) = 0 \end{cases}$$

Ch13. Neural Process

- Want: functional distribution \rightarrow ~~etc.~~

* Gaussian Process (GP)

: Nonparametric Bayesian method

\rightarrow data-efficient (tailored for sparse data)

\rightarrow Inference \rightarrow ~~etc.~~

* Deep neural network (DNN)

\rightarrow data \rightarrow ~~etc.~~

\rightarrow GP-like inference \rightarrow ~~etc.~~

* Neural process

: Neural Network based probabilistic model

* Generative Query Network

: At the viewing θ and dataset \mathcal{D} view query \mathbf{z} \rightarrow generate \mathbf{x}

* Conditional Neural Process (CNP)

: NN that generates θ

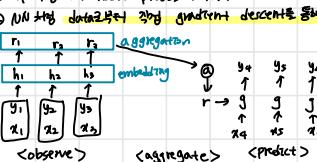
Table: MV GP CNP

Can fit more than 1 function at test time	X	O	O
cheap	O	X	O

① Spatio-temporal stochastic processes \rightarrow ~~etc.~~

② NN that don't need 2nd order gradient descent to find truth

Diagram:



Ex. One-shot Classification

- model: $L_\theta(f(T)|0, T) \approx p_\theta(f(T)|0, T)$

O: context points, T: target inputs

① 1st & 2nd permutations invariant

\rightarrow input의 순서에 관계없이 output은 %

- Architecture

: $\mathbf{r} = h_\theta(\mathbf{x}, \mathbf{y}_n)$, $\mathbf{r} = \alpha(\mathbf{x}_1, \dots, \mathbf{x}_n)$. $\mathbf{r}_n = g_\theta(\mathbf{x}_n, \mathbf{r})$

② regression: $\mathbf{r}_n = g_\theta(\mathbf{x}_n, \mathbf{y}_n)$

③ classification: class probabilities el logits \rightarrow output

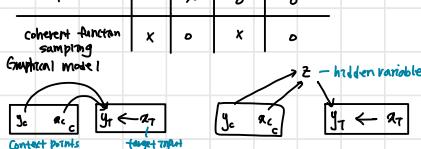
* Neural process

: CNP el %

Can fit more than 1 function at test time	X	O	O
cheap	O	X	O

coherent function sampling	X	O	X
Graphical model!			

- Graphical model!



① encoder: $\mathbf{r}_n = h_\theta(\mathbf{x}_n, \mathbf{y}_n)$

② Aggregator: $\mathbf{r} = \alpha(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \rightarrow$ Gaussian RV el θ 의 M.S.를 parameterize

③ Conditional Decoder: $\mathbf{r}_n = g_\theta(\mathbf{x}_n, \mathbf{z})$

\rightarrow CNP의 추가된 부분