

E-Deep 2회차 발표 (1/14) 논문 리뷰

"Generative Adversarial Nets"

Abstract

이 논문에서는 2가지 모델을 동시에 학습시킨다 : Generative model인 G ,
Discriminative model인 D 이다.
Training 방식은 G 는 D 가 실수를 하도록,
(G 가 만들어낸 sample과 진짜 Data인 training data를 구분하지 못하도록,) 이루어진다. G 와 D 가 마치 minmax 게임을 하는 것처럼 진행된다.
임의의 G, D 에 대해 unique solution이 존재한다.

1. Introduction

G : Generative Model \rightarrow Training Data의 분포를 파악해서 이미지들을 만들어냄.
 D : Discriminative Model \rightarrow G 가 만든 sample과 Data를 구분.

이제껏, (~2014) 딥러닝은 Discriminative Model 계에서 큰 성공을 이루어왔다.
(High-dimensional data를 class Label로 매핑)

이 딥러닝 모델은 최적화과정에서 Dropout과 Backpropagation을 베이스로 한다. 이 최적화는 piecewise linear units을 사용하는데,
Deep Generative models은 piecewise linear units을 사용하기 힘들고,
다중기함수 확률 계산으로 인해 Dropout과 Backpropagation의 영향이 적다.

이 논문에서 제안하는 Adversarial Nets에서는, Generative model이 Discriminative model과 대립하게 된다.

비유> G : 위조 지폐 만드는 범죄자

D : 위조 지폐 탐지하는 경찰.

) 시간이 갈수록 G 는 더 위조지폐를 잘 만들고,
 D 는 더 구분하는 능력 향상

\Rightarrow 나중에는 G 가 만든 위조지폐와 진짜 지폐가 구분하기 힘들 정도로 유사해진다!

(2. 생략)

3. Adversarial Nets

Data x 에 대해 generator의 분포 P_g 을 러닝하기 위해, input 노이즈 변수 $P_z(z)$ 에 대해 prior를 정의한다. 그러면,

그리고 나서 Data space로의 매핑을 $G(z; \theta_g)$ 로 나타낸다.

(G : parameter θ_g 로 표현된 multilayer perceptron으로 나타내진 미분가능한 함수).

또한, second multilayer perceptron $D(x; \theta_d)$ 도 정의한다.

($D(x)$ 는 데이터가 P_g 로부터 오지 않은 x (output: a scalar. i.e. Input이 진짜 Data인지 아닌지) 확률.)

* Training * ① D : trained to maximize probability

(데이터가 G 에서 왔는지 training data에서 왔는지 올바른 라벨을
발하는 확률)

② G : maximize $\log(1 - D(G(z)))$ 즉, $D(G(z))$ 가 커지도록
trained to

($G(z)$ 를 D 가 data label로
표기하도록)

식으로 표현하면,

$$\min_G \max_D V(D, G) = \min_G \max_D \mathbb{E}_{x \sim P_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim P_z(z)} [\log(1 - D(G(z)))]$$

< Algorithm 1 >

k : a hyperparameter

for number of training iterations do

for k steps do

- sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$ 만든 데이터이므로 noise sample로 표현
- sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{data}(x)$. [진짜 Data 들
- Update D (using stochastic gradient) "ascending" ← Do it to Maximize

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))]$$

end for

- sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from $p_g(z)$.
- Update G "descending" ← Go to Minimize

$$\nabla_{\theta_g} \sum_{i=1}^m \log (1 - D(G(z^{(i)})))$$

end for

4. Theoretical Results

$\min_G \max_D V(G, D)$ 가 global optimum 을 가진다. ($P_g = P_{data}$)

4.1. Global optimality of $P_g = P_{data}$.

(먼저 G fixed 상태에서 D 를 K steps 학습시키면)

<Proposition 1> For G fixed, the optimal discriminator D is

$$D_G^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)}.$$

Pf) Given G , to maximize $V(G, D)$,

$$V(G, D) = \mathbb{E}_{x \sim P_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim P_g(z)} [\log (1 - D(G(z)))]$$

$$\begin{aligned} \mathbb{E}_{x \sim P_{data}(x)} [g(x)] &= \int p(x) g(x) dx \\ &= \int_x P_{data}(x) \log D(x) + \int_z P_g(z) \log (1 - D(G(z))) \\ &= \int_x P_{data}(x) \log D(x) + P_g(x) \log (1 - D(x)) dx \end{aligned}$$

Lemma

For $\forall (a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$, the function $y \rightarrow a \log(y) + b \log(1-y)$ 는 $y = \frac{a}{a+b}$ 에서 maximum 을 가진다. in $[0, 1]$.

$$\left(\frac{d}{dy} (a \log y + b \log(1-y)) = \frac{a}{y} + \frac{-b}{1-y} = \frac{a - ay - by}{y(1-y)} = 0 \text{ when } y = \frac{a}{a+b} \right)$$

위에서 a 는 $P_{data}(x)$, b 는 $P_g(x)$ 다. $P_{data}(x)$ 와 $P_g(x)$ 모두 0 이 되지 않으므로
($\sum, \text{Supp}(P_{data}) \cup \text{Supp}(P_g)$ 안에서 정의됨)

참고 $\text{supp}(f) = \overline{\{x \mid f(x) \neq 0\}} \supset \{x \mid f(x) \neq 0\}$.

위의 Lemma에 의해서,

$$V(G, D) = \int_x \underbrace{P_{\text{data}}(x) \log D(x) + P_g(x) \log(1-D(x))}_{D(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \text{에서 Maximum을 가짐}} dx$$

G 는 fixed. Proposition 1에서 D 는 D_G^* 을 가진다. max값으로

$$\begin{aligned} C(G) &= \max_D V(G, D) = \mathbb{E}_{z \sim P_{\text{data}}} [\log D_G^*(z)] + \mathbb{E}_{z \sim P_g} [\log(1 - D_G^*(G(z)))] \\ &= \mathbb{E}_{z \sim P_{\text{data}}} \log \left[\frac{P_{\text{data}}(z)}{P_{\text{data}}(z) + P_g(z)} \right] + \mathbb{E}_{z \sim P_g} \log \left[\frac{P_g(z)}{P_{\text{data}}(z) + P_g(z)} \right] \end{aligned}$$

thm 1

The global minimum of $C(G)$ is achieved

$$\iff P_g = P_{\text{data}}$$

At that point, $C(G) = -\log 4$.

Pf) (\Leftarrow) For $P_g = P_{\text{data}}$, $D_G^* = \frac{P_{\text{data}}}{P_{\text{data}} + P_g} = \frac{1}{2}$.

Hence, $C(G) = \mathbb{E} \left[\log \frac{1}{2} \right] + \mathbb{E} \left[\log \frac{1}{2} \right] = 2 \times \log \frac{1}{2} = -\log 4$

(\Rightarrow)

$$C(G) = V(G, D_G^*)$$

$$\begin{aligned} C(G) - (-\log 4) &= \mathbb{E}_{z \sim P_{\text{data}}} \left[\log \frac{P_{\text{data}}}{P_{\text{data}} + P_g} \right] + \mathbb{E}_{z \sim P_g} \left[\log \frac{P_g}{P_{\text{data}} + P_g} \right] - (-\log 4) \\ &= \mathbb{E}_{z \sim P_{\text{data}}} \left[\log \frac{P_{\text{data}} \times 2}{P_{\text{data}} + P_g} \right] + \mathbb{E}_{z \sim P_g} \left[\log \frac{P_g \times 2}{P_{\text{data}} + P_g} \right] \end{aligned}$$

$\log 2 + \log 2$

$$CCG) - (-\log 4) = \mathbb{E}_{x \sim P_{data}} \left[\log \frac{2P_{data}}{P_{data} + P_g} \right] + \mathbb{E}_{x \sim P_g} \left[\log \frac{2P_g}{P_{data} + P_g} \right]$$

$$= \underbrace{KL \left(P_{data} \parallel \frac{P_{data} + P_g}{2} \right) + KL \left(P_g \parallel \frac{P_{data} + P_g}{2} \right)}_{\substack{\uparrow \\ KL(P \parallel Q) = \mathbb{E}_{x \sim P(x)} \log \frac{P(x)}{Q(x)}}}$$

$$KL(P \parallel Q) = \mathbb{E}_{x \sim P(x)} \log \frac{P(x)}{Q(x)}$$

$$= 2JSD(P_{data} \parallel P_g)$$

$$(JSD(P \parallel Q) = \frac{1}{2} KL(P \parallel Q) + \frac{1}{2} (Q \parallel M))$$

$$\Rightarrow CCG) = -\log 4 + 2JSD(P_{data} \parallel P_g)$$

JSD 는 항상 nonnegative 이고, 두 분포가 같을 때만 0 이다.

$$\text{즉, } CCG) = -\log 4 + 2JSD(P_{data} \parallel P_g) \geq -\log 4$$

즉, CCG) 가 global minimum을

~~CCG)~~ ~~(P_{data} = P_g일 때 등호성립)~~

가져오는 건 $P_{data} = P_g$ 를 뜻한다. ◻



< Proposition 2 >

If G and D have enough capacity, at each step of Algorithm 1, D is allowed to reach its optimum given G and P_g is updated so as to improve the criterion

$$\mathbb{E}_{x \sim P_{data}} [\log D_G^*(x)] + \mathbb{E}_{x \sim P_g} [\log (1 - D_G^*(x))]$$

$$\Rightarrow P_g \rightarrow P_{data}.$$

즉, Algorithm 1의 Iterative 방식은 Generator G 가 Data의 분포물에 수렴함.

Pf) Consider $V(G, D) = U(P_g, D)$

Note $U(P_g, D)$ is convex in P_g . (Global optimum이 존재.)

If $f(x) = \sup_{d \in A} f_d(x)$ and $f_d(x)$ is convex in $x \forall d$,

then $\partial f_\beta(x) \in \partial f$ if $\beta = \operatorname{argsup}_{d \in A} f_d(x)$.

$\sup_D U(P_g, D)$ is convex in P_g with a unique global optima.
($-\log 4$ when $P_g = P_{data}$.)



실제 실험에서는 θ_g 를 업데이트 시킴으로써 P_g 를 optimize. (즉, G 의 parameter를 update)

2. Multilayer perceptron을 사용하면 Critical points가 여러개 나온다. (Min, Max, saddle points 등 AA)

그래도 성능은 잘 나온다!

5. Experiments

G는 ReLU, Stigmod 를 섞어서 사용한다. D는 maxout 사용.

또한, D를 학습시키는 데에 Dropout 이 사용된다.

6. Advantages and disadvantages

단점 : • No explicit representation of $P_g(x)$.

- D가 학습중에 G와 잘 동기화되어야함. 특히, D를 업데이트하지 않고 G를 과하게 학습시키면 안됨.

장점 : • Markov Chain이 필요없다. Backprop만 사용!

- 사용가능한 함수들 다
- No Inference is needed.
- G가 Data로부터 바로 학습되는 것이 아니라 D에서 계산한 gradients 이용.
- Sharp, degenerate 분포를 나타낼 수 있다.
(MC 에서는 흐릿함).