

type

연구

AI summary

Try on this page

[이과대학 수학] 알핀식 Nov 25

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio
GAN 구조 학습

RESEARCH

OVERVIEW

generative model G: 이미지 생성

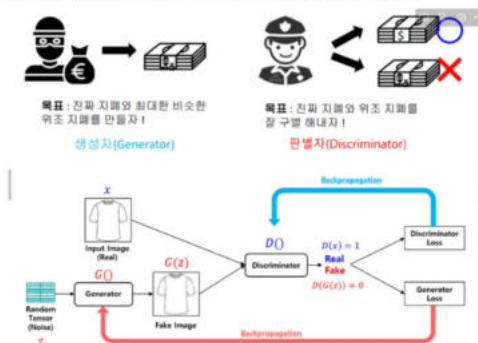
discriminative model D: 이미지 구별

G는 D의 실수를 최대화 시키는게 목표이다. minimax tow-player game과 동일하다.

CONTENT

1. Introduction

딥러닝은 class label을 mapping하는 모델에서 Backpropagation과 Dropout 등을 사용하여 뛰어난 성과를 거두었다. 하지만 deep generative model에서는 MLE에서 intractable 확률 계산을 approximating하는 문제에 어려움을 겪어왔다. 그래서 저자는 새로운 generative model을 제시해서 이러한 어려움을 피해나갈 수 있었다.



adversarial nets에서 discriminative model은 샘플이 model distribution으로부터 왔는지 아닌지 확인하고 generative model은 fake 이미지를 생성해서 속이는 역할을 한다.

논문에서 저자는 generative model을 counterfeiters(가짜 위조자), discriminative model을 경찰에 비유했다.

generative model은 랜덤 노이즈 Z를 multilayer perceptron에 넣어 fake image를 만들고 discriminative model은 multilayer perceptron에 데이터를 넣어 모방 데이터인지 진짜 데이터인지 확인한다. generative model과 discriminative model 2개를 합쳐 adversarial nets라고 한다. discriminative loss와 generator loss 2개를 뒤서 각각에 대해 Backpropagation을 사용하고 dropout algorithms를 사용하였다. 여기서 중요한 것은 approximate inference나 Markov chain을 사용하지 않았다는 점이다.

2. Adversarial nets

 p_g : generator's distribution over data x
 $p_z(z)$: input noise variables

 $G(z; \theta_g)$: mapping to data space with parameters θ_g
 $D(x; \theta_d)$: represents probability that x came from the data rather than p_g = single scalar

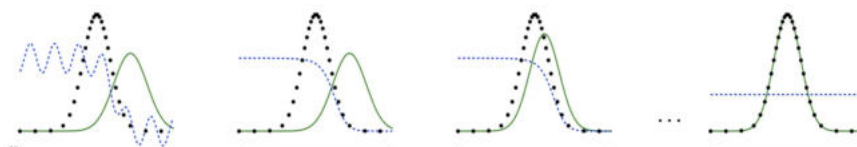
< Train >

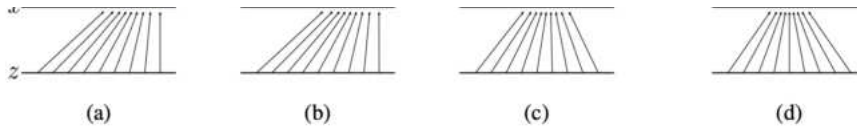
 D : maximize the probability of assigning the correct label to both training examples and samples from G
 G : minimize $\log(1 - D(G(z)))$

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

- Discriminator는 진짜 데이터는 1로 판단하여 첫 번째 항을 0으로, 두 번째 가짜 데이터를 0으로 판단하여 두 번째 항도 0으로 만드는 것이 목표이다.
- 즉, Discriminator가 얻을 수 있는 이상적인 최대값은 0이다.
- 반대로 Generator는 두 번째 항의 $D(G(z)) = 1$ 로 판단하게 하여 두 번째 항 전체를 $-\infty$ 만드는 것이 목표이다.
- 즉, Generator가 얻을 수 있는 이상적인 최솟값은 $-\infty$ 이다.
- 학습을 계속하다보면 Discriminator도 어느 것이 진짜인지 가짜인지 구분하기 힘들어져서 $D(x) = 0.5$ 에 수렴하게 된다.

2.1 Theoretical analysis of adversarial nets





검은색: p_{data} , 파란색: discriminative distribution, 초록색: generative distribution

- 처음에는 원본 데이터에 대해서는 1에 가까운 값을 fake 이미지에 대해서는 0에 가까운 값으로 구별을 한다.

$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

- D의 결과값이 위의 수식으로 수렴한다. 추후 증명예정.
- 그러다가 generative distribution이 점점 원본 데이터의 분포와 같아지게 되고 가짜를 진짜와 가짜를 구별해내지 못하게 되어 0.5로 수렴한다.
- 결론적으로 $p_{data} = p_g$ 일 때 global optimum을 가진다.

2.2 Theoretical Results

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

```

for number of training iterations do ..... 1
  for  $k$  steps do ..... 2
    • Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ . ..... 3
    • Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution ..... 4
       $p_{data}(x)$ .
    • Update the discriminator by ascending its stochastic gradient:
      
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$
 ..... 5
    end for
    • Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
    • Update the generator by descending its stochastic gradient:
      
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$
 ..... 6
  end for
  The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.
  
```

- epoch만큼 반복문 수행
- k번만큼 discriminator gradient ascending 시행(max) → 자라는 k=1로 됨
- 1번만큼 generator gradient descend 시행(min) → 실제로는 gradient ascend 사용

2.3 Global Optimality

Proposition 1) G와 D를 정하고, discriminator p_g

$$D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

$$P_G^*(V(G, D)) = \int_{\mathcal{X}} p_{\text{data}}(x) \log D_G^*(x) dx + \int_{\mathcal{Z}} p_g(z) (1 - D_G^*(G(z))) dz$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) \log D_G^*(x) + p_g(x) (1 - \log D_G^*(x)) dx$$

For $\forall (a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$ $y \mapsto a \log y + b \log(1 - y)$

$$\Rightarrow \text{a.s. } \frac{a}{y} - \frac{b}{1-y} = 0, \quad a(1-y) - by = 0, \quad y(a+b) = a, \quad y = \frac{a}{a+b}$$

$$\therefore D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

$$C(G) = \max_D V(G, D)$$

$$= E_{x \sim p_{\text{data}}} [\log D_G^*(x)] + E_{z \sim p_g} [\log(1 - D_G^*(G(z)))]$$

$$= E_{x \sim p_{\text{data}}} [\log D_G^*(x)] + E_{x \sim p_g} [\log(1 - D_G^*(x))]$$

$$= E_{x \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right] + E_{x \sim p_g} \left[\log \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right]$$

Theorem 1) The global minimum of the Variational training criterion $C(G)$ is achieved if and only if $p_g = p_{\text{data}}$. At that point $C(G)$ achieves the value $-\log 4$. *Gloptimality, C(G)에 최소값*

Proof) $C(G) \leq -\log 4$ $p_g = p_{\text{data}} \Rightarrow D_G^*(x) = \frac{1}{2}$

$$C(G) = \log 2 + \log 2 = -\log 4$$

$$D_{KL}(p \| q) = \int p(x) \log \frac{p(x)}{q(x)}$$

$$JSD(p \| q) = \frac{1}{2} D_{KL}(p \| M) + \frac{1}{2} D_{KL}(q \| M) \geq 0, \quad M = \frac{p+q}{2}$$

$$C(G) = C(G) + \log 4 - \log 4$$

$$= -\log 4 + E_{x \sim p_{\text{data}}} \left[\log \frac{2 \times p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right] + E_{x \sim p_g} \left[\log \frac{2 \times p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right]$$

$$= -\log 4 + KL(p_{\text{data}} \| \frac{p_{\text{data}} + p_g}{2}) + KL(p_g \| \frac{p_{\text{data}} + p_g}{2})$$

$$= -\log 4 + 2 \times JSD(p_{\text{data}} \| p_g) \geq -\log 4$$

Where $p_{\text{data}} = p_g$

- 2개의 증명들 통해 minmax problem이 global minimum에서 unique solution을 가지고 어떠한 조건이 만족하면 그 solution이 값으로 수렴한다.

2.4 Convergence of Algorithm 1

Proposition 2. If G and D have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given G , and p_g is updated so as to improve the criterion

$$E_{x \sim p_{\text{data}}} [\log D_G^*(x)] + E_{x \sim p_g} [\log(1 - D_G^*(x))]$$

then p_g converges to p_{data}

Proof. Consider $V(G, D) = U(p_g, D)$ as a function of p_g as done in the above criterion. Note that $U(p_g, D)$ is convex in p_g . The subderivatives of a supremum of convex functions include the derivative of the function at the point where the maximum is attained. In other words, if $f(x) = \sup_{\alpha \in A} f_\alpha(x)$ and $f_\alpha(x)$ is convex in x for every α , then $\partial f_\beta(x) \in \partial f$ if $\beta = \arg \sup_{\alpha \in A} f_\alpha(x)$. This is equivalent to computing a gradient descent update for p_g at the optimal D given the corresponding G . $\sup_D U(p_g, D)$ is convex in p_g with a unique global optima as proven in Thm 1, therefore with sufficiently small updates of p_g , p_g converges to p_{data} , concluding the proof. ■

In practice, adversarial nets represent a limited family of p_g distributions via the function $G(z; \theta_g)$, and we optimize θ_g rather than p_g itself. Using a multilayer perceptron to define G introduces multiple critical points in parameter space. However, the excellent performance of multilayer perceptrons in practice suggests that they are a reasonable model to use despite their lack of theoretical guarantees.

- 문제 2: G, D가 충분한 용량을 갖고 Algorithm 1 각 step에서 discriminator는 주어진 G에 대하여 위에서 제시한 목적의 값으로 도달하고 p_g를 업데이트 해서 criterion을 개선하면 p_g → p_data로 수렴한다.
- 정리하자면 주어진 Generator에 대해 D가 optimum으로 수렴하고, p_g를 업데이트하면 p_g → p_data로 수렴한다는 이야기이다.

$$V(G, D) = U(p_g, D) : p_g \text{의 함수로 생각}$$

$$U(p_g, D) : \text{convex in } p_g$$

볼록 함수에서 최상부의 도함수 = 최대값에 도달한 지점의 도함수

$$\text{if } f(x) = \sup_{\alpha \in A} f_\alpha(x) \text{ and } f_\alpha(x) \text{ is convex in } x \text{ for every } \alpha, \text{ then } \partial f_\beta(x) \in \partial f \text{ if } \beta = \arg \sup_{\alpha \in A} f_\alpha(x)$$

- 정확하게 이해가 가지는 않는다.
- p_g에 대한 gradient descent update와 동일하다.
- convex 하므로 global optimum을 가지고 p_g는 결국 p_data로 수렴하게 된다.
- Subderivative 정의

수학에서 하방미분(subdifferential, subderivative)은 미분을 일반화하여 미분가능하지 않은 볼록 함수에 적용할 수 있도록 하는 방법이다. 볼록 최적화 등 볼록 함수를 연구하는 해석에서 중요하게 사용된다.

정의 (한글)

볼록함수 $f: I \rightarrow \mathbb{R}$ 가 있을 때, I 의 점 x_0 에서의 하방미분계수는:

$$f(x) - f(x_0) \geq c(x - x_0)$$

가 I 의 모든 점 x 에 대해 성립하게 하는 실수 c 를 가리킨다.

3. Experiments

- G: ReLU(Rectifier linear activations) + sigmoid 은합 사용
- D: maxout 사용, 학습시일 때 dropout 사용



- 장점: backpropagation만 사용, inference 필요없음, 이미지가 더 선명함
- 단점: $p_g(x)$ 가 명시적으로 존재하지 않음, D 와 G 가 균형을 잘 맞추지 못함 (학습되어야 할, 만약 G 가 너무 발전하기 이전에 발전해버리면 G 가 z 데이터를 너무 많이 불러오기 때문이다).

단어장

adversarial: 적대적인

pitted: 싸우게 하다

counterfeits: 화폐 위조자

SEARCH

REFER

#	Title	type	link
(1)			

Proposition 1) G 가 고정된 경우, $\exists!$ 의 discriminator D ,

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$Pf) V(G, D) = \int_x p_{data}(x) \log D(x) dx + \int_z p_z(z) (1 - D(G(z))) dz$$

$$= \int_x [p_{data}(x) \log D(x) + p_g(x) (1 - \log D(x))] dx$$

$$\text{For } \forall (a, b) \in \mathbb{R}^2 \setminus (0, 0) \quad y \mapsto a \log y + b \log(1-y)$$

$$\Rightarrow \text{미분: } \frac{a}{y} - \frac{b}{1-y} = 0, \quad a(1-y) - by = 0, \quad y(a+b) = a, \quad \boxed{y = \frac{a}{a+b}}$$

$$\therefore D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$C(G) = \max_D V(G, D)$$

$$= E_{x \sim p_{data}} [\log D_G^*(x)] + E_{z \sim p_z} [\log (1 - D_G^*(G(z)))]$$

$$= E_{x \sim p_{data}} [\log D_G^*(x)] + E_{x \sim p_g} [\log (1 - D_G^*(x))]$$

$$= E_{x \sim p_{data}} \left[\log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \right] + E_{x \sim p_g} \left[\log \frac{p_g(x)}{p_{data}(x) + p_g(x)} \right]$$