

<한생배경>

Reviewer : 김정학

다음의 residual network 를 살펴보자.

$$h_{t+1} = h_t + f(h_t, \theta_t)$$

이 식을 간단하게 되면 Euler method와 상당히 유사한 것을 확인할 수 있다.

만약 우리가 hidden layer를 무수히 많이 추가하면 어떻게 될까?

그러면 위 식의 $f(h_t, \theta_t)$ 를 h 의 변화량으로 볼 수도 있다. 즉,

$$\frac{dh(t)}{dt} = f(h(t), t, \theta) \text{로 볼 수 있는 것이다.}$$

이는 만약 우리가 $h(T)$ 를 알고 싶으면 $h(0)$ 에 '더하기'만 반복
계속 해서 $h(T)$ 를 얻을 수 있다는 것을 의미한다.

<특징>

① memory efficiency

기존의 backtracking이 아닌 방식으로 함수를 진행하게 되는데
이 방식은 또 다른 덧셈의 집합이라 forward에 대한 정보를
계속 가질 필요가 없어 메모리 관점에서 효율적이다.

↙ 참고

https://leonardoraujosantos.gitbook.io/artificial-intelligence/machine_learning/supervised_learning/backpropagation

(기존 램웨어 메모리 사용이 많은 양이)

② adaptive computation

forward 혹은 backward를 정보로 제공할 것이다. 그에 정보를
제공해 최솟값을 어떻게 할까? 매우 작은 수로 과정을 하는 것은 바람직
이다. 그래서 정보대량(함수)의 변화에 따라서 step-size를 정하는
(adaptive) 방식의 ODE solver를 사용한다 (e.g dopri)

③ Scalable and invertible normalizing flows

데이터가 function을 통해서 변환될 때 분포의 변화를

계산하는 과정이 간단해진다. (정확히는 참...)

④ Continuous time-series models

최대 배치가 연속적으로 주어질 수 있어서 시계열 모델에 적용이 될 수 있다.

< backward >

forward는 단순한 직분임을 알았다. 이 논문에서 주를 다루는 내용은 backward를 어떻게 하는지이다.

우선 우리는 다음 식을 통해서 ODE layer를 통과한 결과를 알 수 있다.

$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f(z(t), t, \theta) dt$$

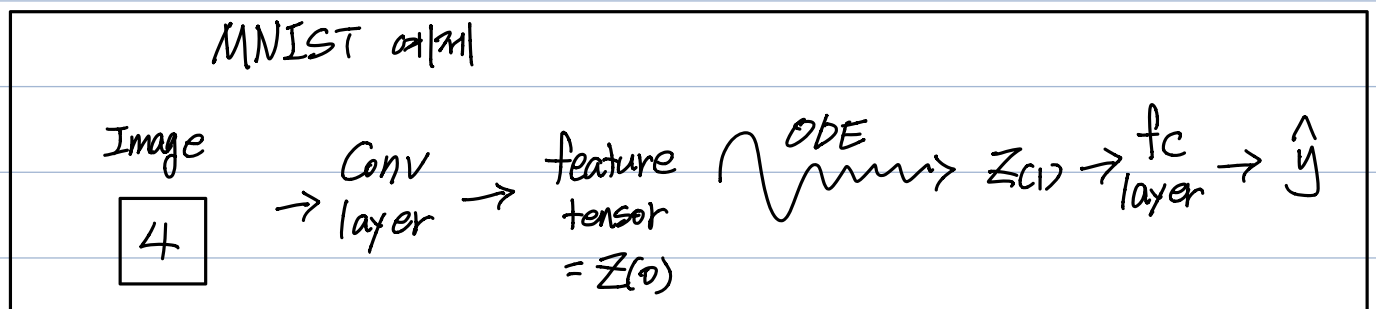
$$= \underset{\text{적분기}}{\text{ODE Solver}} \left(\underset{\text{초기값}}{z(t_0)}, \underset{\text{ODE Func}}{f}, \underset{\text{시작시간}}{t_0}, \underset{\text{종료시간}}{t_1}, \underset{\text{parameters}}{\theta} \right)$$

이로 Loss를 주하면

$$L(z(t_1)) = L \left(z(t_0) + \int_{t_0}^{t_1} f(z(t), t, \theta) dt \right)$$

$$= L \left(\text{ODE Solver} \left(z(t_0), f, t_0, t_1, \theta \right) \right)$$

여기서 우리가 학습하기 위해 필요한 미분값을 보자
우선 ODE network는 다음과 같다.



우리는 $z(t_0)$ 을 뽑아주는, 즉 feature를 뽑아주는 network (여기에서는 Conv layer)를 학습하기 위해

- ① $\frac{dL}{dz(t_0)}$ 가 필요하다. 그리고 ODE layer를 학습하기 위해
- ② $\frac{dL}{d\theta}$ 가 필요하다.

이를 구하기 위해 adjoint $a(t) = \frac{dL}{dz(t)}$ 를 도입한다.

이를 수식으로 아래처럼 해보면 (논문 Appendix를 참고)

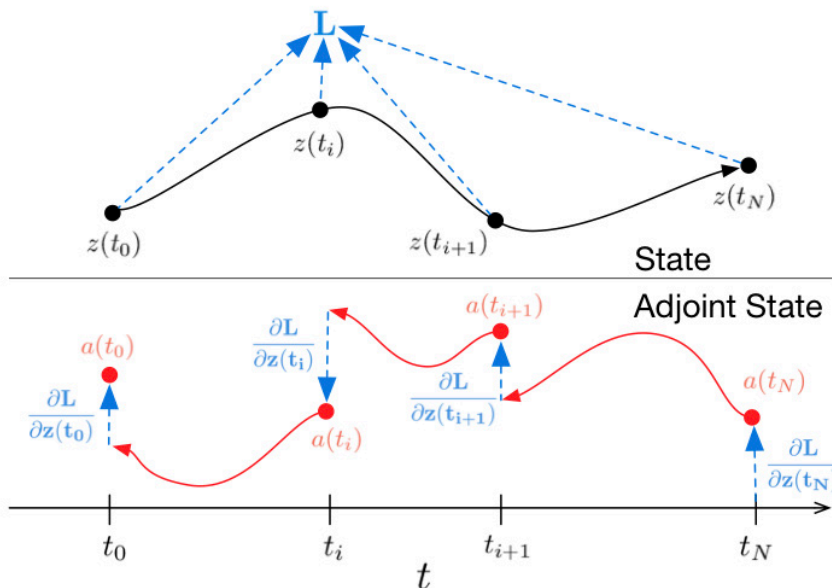
$$\frac{da(t)}{dt} = -a(t)^T \frac{df(z(t), t; \theta)}{dz}$$

를 얻을 수 있다.

가만 살펴보자 어차피 ① $a(t_0) = \frac{dL}{dz(t_0)}$ 를 구하기 위해서

$\frac{da(t)}{dt}$ 를 이용해 똑같은 ODE Solver를 계산하면 된다.

그리고 또 만약에 ② $\frac{dL}{d\theta} = - \int_{t_1}^{t_0} a(t)^T \frac{df(z(t), t; \theta)}{d\theta} dt$ 를 얻을 수 있다.



이런 단계를 나열한 도식이다.

- State에서는 평상 ODE Solver를 계산한다.
- adjoint state에서는 augmented dynamic을 이용한다.
시간점 사이사이에서는 ODE Solver를 계산하고 각 시간점에서
업데이트를 한다.

이것이 backpropagation을 위한 알고리즘이다.

초기 상태

Algorithm 1 Reverse-mode derivative of an ODE initial value problem

Input: dynamics parameters θ , start time t_0 , stop time t_1 , final state $\mathbf{z}(t_1)$, loss gradient $\partial L / \partial \mathbf{z}(t_1)$
 $s_0 = [\mathbf{z}(t_1), \frac{\partial L}{\partial \mathbf{z}(t_1)}, \mathbf{0}_{|\theta|}]$ \triangleright Define initial augmented state
def aug_dynamics($[\mathbf{z}(t), \mathbf{a}(t), \cdot], t, \theta$): \triangleright Define dynamics on augmented state
 return $[f(\mathbf{z}(t), t, \theta), -\mathbf{a}(t)^\top \frac{\partial f}{\partial \mathbf{z}}, -\mathbf{a}(t)^\top \frac{\partial f}{\partial \theta}]$ \triangleright Compute vector-Jacobian products
 $[\mathbf{z}(t_0), \frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}] = \text{ODESolve}(s_0, \text{aug_dynamics}, t_1, t_0, \theta)$ \triangleright Solve reverse-time ODE
return $[\frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}]$ \triangleright Return gradients

feature를 뽑아내는
network 층

ODE layer 층.

반대방향을 걸음.