September 1, 2020



# **Continuous Normalizing Flows**

# Vikram Voleti PhD student, Mila, University of Montreal <u>Supervisor</u>: Prof. Christopher Pal

voletiv.github.io

@ (virtual) Mila, Montreal, Canada



# **1. Ordinary Differential Equations**

- 2. Neural ODEs
- 3. Continuous Normalizing Flows
- 4. Later research



$$rac{dx(t)}{dt}=f(x(t),t, heta); \;\; x(t_0)$$
 is given;  $\; x(t_1)=\; ?$ 

Solution:

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

$$egin{aligned} & \displaystyle rac{ ext{Example:}}{rac{dx}{dt}=2t; \; x(0)=2; \; x(1)=\;? \ \Rightarrow x(1)=x(0)+\int_{0}^{1}2t\; dt \ &=x(0)+(t^{2}|_{t=1}-t^{2}|_{t=0}) \ &=2+1^{2}-0^{2} \ &=3 \ \end{aligned}$$

### Vikram Voleti

$$rac{dx(t)}{dt}=f(x(t),t, heta); \;\; x(t_0)$$
 is given;  $\;x(t_1)=\;?$ 

Solution: $x(t_1)=x(t_0)+\int_{t_0}^{t_1}f(x(t),t, heta)~dt$ 

What if this cannot be analytically integrated?

$$egin{aligned} rac{\mathrm{Example:}}{rac{dx}{dt} &= 2xt \ ; \ x(0) = 3 \ \Rightarrow \int rac{1}{2x} \ dx &= \int t \ dt \ \Rightarrow rac{1}{2} \mathrm{log} \ x &= rac{1}{2} t^2 + c_0 \ \Rightarrow x(t) &= c e^{t^2} \ x(0) = 3 \Rightarrow c = 2 \ \therefore x(t) = 2 e^{t^2} \ \Rightarrow x(1) = 5.436 \end{aligned}$$

### Vikram Voleti





$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given}; \ \ x(t_1)=\ ?$$

Solution:

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

Approximations to  $\int_{t_0}^{t_1} f(x(t),t, heta) \ dt$ 

# i.e. Numerical Integration :

- Euler method
- Runge-Kutta methods

•

•••



$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given; } \ x(t_1)=\ ?$$

Solution:

 $f(x(t_1) = x(t_0) + \int_{t_0}^{t_1} f(x(t),t, heta) \; dt$ 

1st-order Runge-Kutta / Euler's method:

$$t_{n+1} = t_n + h$$
 ----- Step size  $h$   
 $x(t_{n+1}) = x(t_n) + hf(x(t_n), t_n)$  -  $\blacktriangleright$  Update using derivative



#### Vikram Voleti



$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given}; \ \ x(t_1)=\ ?$$

Solution:

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

## 1st-order Runge-Kutta / Euler's method:

$$egin{aligned} t_{n+1} &= t_n + h \ x(t_{n+1}) &= x(t_n) + hf(x(t_n),t_n) \end{aligned}$$

Example:  $\frac{dx}{dt} = f(x,t) = 2xt ; x(0) = 3; x(1) = ?$ (Solution:  $x(t) = 2e^{t^2}; x(1) \neq 5.436$ )

$$\begin{array}{l} h = 0.25 \\ x(0.25) = x(0) + 0.25 * f(x(0), 0) \\ = 3 + 0.25 * (2 * 3 * 0) \\ = 3 \\ x(0.5) = x(0.25) + 0.25 * f(x(0.25), 0.25) \\ = 3 + 0.25 * (2 * 3 * 0.25) \\ = 3.375 \\ x(0.75) = x(0.5) + 0.25 * f(x(0.5), 0.5) \\ = 3.375 + 0.25 * (2 * 3.375 * 0.5) \\ = 4.21875 \\ x(1) = x(0.75) + 0.25 * f(x(0.75), 0.75) \\ = 4.21875 + 0.25 * (2 * 4.21875 * 0.75) \\ = 5.8008 \end{array}$$

### Vikram Voleti



$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given}; \ \ x(t_1)=\ ?$$

Solution:

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

## 1st-order Runge-Kutta / Euler's method:

$$egin{aligned} t_{n+1} &= t_n + h \ x(t_{n+1}) &= x(t_n) + hf(x(t_n),t_n) \end{aligned}$$



### Vikram Voleti



$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given}; \ \ x(t_1)=\ ?$$

Solution:

$$\int x(t_1) = x(t_0) + \int_{t_0}^{t_1} f(x(t),t, heta) \ dt$$

---- Default ODE solver used in MATLAB: <u>https://blogs.mathworks.com/loren/2015/09/23/o</u> <u>de-solver-selection-in-matlab/</u>

### Vikram Voleti



$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given}; \ \ x(t_1)=\ ?$$

### Solution:

$$\int x(t_1) = x(t_0) + \int_{t_0}^{t_1} f(x(t),t, heta) \ dt$$





$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0)$$
 is given;  $\ x(t_1)=\ ?$ 

### Solution:

 $x(t_1) = ext{ODESolve}(\; f(x(t),t, heta),\; x(t_0),\; t_0,\; t_1\;)$ 



### Vikram Voleti



# 1. Ordinary Differential Equations

# 2. Neural ODEs



$$rac{dx(t)}{dt}=f(x(t),t, heta); \ \ x(t_0) ext{ is given}; \ \ x(t_1)=\ ?$$

Solution:

 $x(t_1) = ext{ODESolve}(\; f(x(t),t, heta),\; x(t_0),\; t_0,\; t_1\;)$ 

f is a neural network!

**Paradigm shift**: whereas earlier *f* was pre-defined/hand-designed according to the domain, here we would like to estimate an *f* that suits our objective.



## ODEs

$$\begin{array}{c} \frac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}(t),t,\theta) \\ \hline \\ \text{Vector} \\ \text{notation} \\ \mathbf{x}_{n+1} = \mathbf{x}_n + h \; f(\mathbf{x}_n,t_n,\theta) \end{array}$$

# Residual networks

$$\mathbf{x}_{l+1} = \operatorname{ResBlock}(\mathbf{x}_l, \theta)$$
  
 $\mathbf{x}_{l+1} = \mathbf{x}_l + g(\mathbf{x}_l, \theta)$   
Skip connection



# ODEs

$$rac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}(t),t, heta)$$
 $igsquare$  Euler discretization
 $\mathbf{x}_{n+1} = \mathbf{x}_n + h \; f(\mathbf{x}_n,t_n, heta)$ 

Forward propagation:  $\mathbf{x}(t_1) = \text{ODESolve}(\ f(\mathbf{x}(t), t, \theta), \ \mathbf{x}(t_0), \ t_0, \ t_1 \ )$ 

 $L(\mathbf{x}(t_1)) o rac{\partial L}{\partial heta}$  How to compute this? Update heta to reduce L

# Residual networks

$$\mathbf{x}_{l+1} = \operatorname{ResBlock}(\mathbf{x}_l, \theta)$$
  
 $\mathbf{x}_{l+1} = \mathbf{x}_l + g(\mathbf{x}_l, \theta)$   
Skip connection

$$\mathbf{y}_{pred} = \operatorname{ResNet}(\mathbf{x})$$

$$\overset{\checkmark}{\blacktriangleright} Stacked \operatorname{ResBlocks}$$

$$L(\mathbf{y}_{pred}) 
ightarrow rac{\partial L}{\partial heta}$$
Update  $heta$ to reduce  $L$ 

https://arxiv.org/pdf/1512.03385.pdf

https://arxiv.org/pdf/1806.07366.pdf

Vikram Voleti



## ODEs

$$rac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}(t),t, heta)$$
 $igsquare$  Euler discretization
 $\mathbf{x}_{n+1} = \mathbf{x}_n + h \; f(\mathbf{x}_n,t_n, heta)$ 

Forward propagation:

$$\mathbf{x}(t_1) = ext{ODESolve}(\ f(\mathbf{x}(t), t, heta), \ \mathbf{x}(t_0), \ t_0, \ t_1)$$

 Back-propagate through the ODE Solver!

High memory cost -

need to save all activations of all iterations of ODESolve.

Can we do better?

Yes.

 $L(\mathbf{x}(t_1)) 
ightarrow rac{\partial L}{\partial heta}$  Update hetato reduce L



 $L( ext{ODESolve}(\ f(\mathbf{x}(t),t, heta),\ \mathbf{x}(t_0),\ t_0,\ t_1\ )) o rac{\partial L}{\partial heta})$ 

# Adjoint method (Pontryagin et al., 1962)

adjoint 
$$\mathbf{a}(t) = rac{\partial L}{\partial \mathbf{x}}$$
;  $rac{d\mathbf{a}}{dt} = -\mathbf{a}(t)^{ op} rac{\partial f(\mathbf{x}(t), t, heta)}{\partial \mathbf{x}}$ 

Forward propagation:  $\mathbf{x}(t_1) = \text{ODESolve}(f(\mathbf{x}(t), t, \theta), \mathbf{x}(t_0), t_0, t_1) \Rightarrow \mathbf{a}(t_1) = \frac{\partial L}{\partial \mathbf{x}(t_1)}$ 

$$rac{\partial L}{\partial heta} = \int_{t_1}^{t_0} - \mathbf{a}(t)^ op rac{\partial f(\mathbf{x}(t),\,t,\, heta)}{\partial heta} \; dt$$



 $L(\text{ODESolve}(\ f(\mathbf{x}(t), t, \theta), \ \mathbf{x}(t_0), \ t_0, \ t_1 \ )) 
ightarrow rac{\partial L}{\partial heta}$ 

# Adjoint method (Pontryagin et al., 1962)

adjoint 
$$\mathbf{a}(t) = rac{\partial L}{\partial \mathbf{x}}$$
;  $rac{d\mathbf{a}}{dt} = -\mathbf{a}(t)^{ op} rac{\partial f(\mathbf{x}(t), t, heta)}{\partial \mathbf{x}}$ 

Forward propagation:  $\mathbf{x}(t_1) = \text{ODESolve}(f(\mathbf{x}(t), t, \theta), \mathbf{x}(t_0), t_0, t_1) \Rightarrow \mathbf{a}(t_1) = \frac{\partial L}{\partial \mathbf{x}(t_1)}$ 

Back-propagation:

 $x(t_0) = \text{ODESolve}(-f(\mathbf{x}(t), t, \theta)), \mathbf{x}(t_1), t_1, t_0)$ 

$$\Rightarrow \mathbf{a}(t_0) = \frac{\partial L}{\partial \mathbf{x}(t_0)} = \text{ ODESolve}(-\mathbf{a}(t)^\top \frac{\partial f(\mathbf{x}(t), t, \theta)}{\partial \mathbf{x}}, \frac{\partial L}{\partial \mathbf{x}(t_1)}, t_1, t_0)$$

$$:: \frac{\partial L}{\partial \theta} = \int_{t_1}^{t_0} -\mathbf{a}(t)^\top \frac{\partial f(\mathbf{x}(t), t, \theta)}{\partial \theta} dt = \text{ODESolve}(-\mathbf{a}(t)^\top \frac{\partial f(\mathbf{x}(t), t, \theta)}{\partial \theta}, \mathbf{0}_{|\theta|}, t_1, t_0)$$



### Forward propagation:

$$\mathbf{x}(t_1) = \text{ODESolve}(f(\mathbf{x}(t), t, \theta), \mathbf{x}(t_0), t_0, t_1)$$
  
Compute  $L(\mathbf{x}(t_1))$ .  
 $\mathbf{a}(t_1) = \frac{\partial L}{\partial \mathbf{x}(t_1)}$ 

**Back-propagation:** 

$$\begin{array}{c} \mathbf{x}(t_{0}) \\ \frac{\partial L}{\partial \mathbf{x}(t_{0})} \\ \frac{\partial L}{\partial \theta} \end{array} = \text{ODESolve} \left( \begin{array}{c} f(\mathbf{x}(t), t, \theta) \\ -\mathbf{a}(t)^{\top} \frac{\partial f(\mathbf{x}(t), t, \theta)}{\partial \mathbf{x}} \\ -\mathbf{a}(t)^{\top} \frac{\partial f(\mathbf{x}(t), t, \theta)}{\partial \theta} \end{array} \right), \quad \begin{array}{c} \mathbf{x}(t_{1}) \\ \frac{\partial L}{\partial \mathbf{x}(t_{1})} \\ \mathbf{0}_{|\theta|} \end{array} \right), \quad t_{1}, t_{0} \right)$$

Update heta to reduce L



https://arxiv.org/pdf/1806.07366.pdf

$$\mathbf{x}(t_0) \longrightarrow f(\mathbf{x}(t), t, \theta) \longrightarrow \mathbf{x}(t_1)$$
  
Neural ODE



https://openreview.net/pdf?id=B1e9Y2NYvS

Neural ODEs describe a homeomorphism (flow).

- They preserve dimensionality.
- They form non-intersecting trajectories.

#### Vikram Voleti



https://arxiv.org/pdf/1806.07366.pdf

$$\mathbf{x}(t_0) \longrightarrow f(\mathbf{x}(t), t, \theta)$$
  
Neural ODE

$$\mathbf{x}(t_0) \triangleleft \mathbf{f}(\mathbf{x}(t), t, \theta)$$
  
Neural ODE

## Neural ODEs are **reversible** models! Just integrate forward/backward in time.





- 1. Ordinary Differential Equations
- 2. Neural ODEs





Train f to maximize the likelihood of the samples from target distribution  $\log p(\mathbf{x})$ 

https://arxiv.org/abs/1810.01367

Vikram Voleti

# **Continuous Normalizing Flows**



**"FFJORD"** 

### (Free-Form Jacobian Of Reversible Dynamics)

Target distribution

Noise distribution



# **Continuous Normalizing Flows**

# FFJORD (ICLR 2019)

Change of variables:

 $\log p(\mathbf{x}_{im}) - \log p(\mathbf{z}) = \log \det |rac{\mathrm{d} f_ heta}{\mathrm{d} \mathbf{x}(t)}|$ 

Instantaneous change of variables:

$$\frac{\partial \log p(\mathbf{x}(t))}{\partial t} = -\mathrm{Tr}\left(\frac{\partial f_{\theta}}{\partial \mathbf{x}(t)}\right)$$

Initial value:  $\begin{bmatrix} \mathbf{x}_{im} \\ 0 \end{bmatrix}$  $\begin{bmatrix} \mathbf{z} \\ \log p(\mathbf{x}_{im}) - \log p(\mathbf{z}) \end{bmatrix} = \int_{t_0}^{t_1} \begin{bmatrix} f_{\theta}(\mathbf{x}(t), t) \\ -\operatorname{Tr}\left(\frac{\partial f_{\theta}}{\partial \mathbf{x}(t)}\right) \end{bmatrix} \mathrm{d}t$ 

Hutchinson's trace estimator:

$$\log p(\mathbf{z}(t_1)) = \log p(\mathbf{z}(t_0)) - \int_{t_0}^{t_1} \operatorname{Tr}\left(\frac{\partial f}{\partial \mathbf{z}(t)}\right) dt$$
$$= \log p(\mathbf{z}(t_0)) - \int_{t_0}^{t_1} \mathbb{E}_{p(\boldsymbol{\epsilon})} \left[\boldsymbol{\epsilon}^T \frac{\partial f}{\partial \mathbf{z}(t)} \boldsymbol{\epsilon}\right] dt$$
$$= \log p(\mathbf{z}(t_0)) - \mathbb{E}_{p(\boldsymbol{\epsilon})} \left[\int_{t_0}^{t_1} \boldsymbol{\epsilon}^T \frac{\partial f}{\partial \mathbf{z}(t)} \boldsymbol{\epsilon} dt\right],$$

https://arxiv.org/abs/1810.01367

### Vikram Voleti



# **Continuous Normalizing Flows**



# FFJORD (ICLR 2019)

Change of variables:

 $\log p(\mathbf{x}_{im}) - \log p(\mathbf{z}) = \log \det |rac{\mathrm{d} f_ heta}{\mathrm{d} \mathbf{x}(t)}|$ 

Instantaneous change of variables:

 $\frac{\partial \log p(\mathbf{x}(t))}{\partial t} = -\mathrm{Tr}\left(\frac{\partial f_{\theta}}{\partial \mathbf{x}(t)}\right)$ 

$$\begin{bmatrix} \mathbf{x}_{im} \\ 0 \end{bmatrix} \\ \begin{bmatrix} \mathbf{z} \\ \log p(\mathbf{x}_{im}) - \log p(\mathbf{z}) \end{bmatrix} = \int_{t_0}^{t_1} \begin{bmatrix} f_{\theta}(\mathbf{x}(t), t) \\ -\text{Tr}\left(\frac{\partial f_{\theta}}{\partial \mathbf{x}(t)}\right) \end{bmatrix} dt$$

CIFAR10		ImageNet64	
BPD	Time	BPD	Time
3.40	≥5 days	-	-

https://arxiv.org/abs/1810.01367

# How to Train your Neural ODE (ICML 2020)

Introduces 2 regularization terms: 1) Kinetic energy of flow 2) Jacobian norm of flow

$$egin{aligned} \mathcal{K}( heta) &= \int_{t_0}^{t_1} \|f(\mathbf{x}(t),t, heta)\|_2^2 \,\mathrm{d}t \ \mathcal{B}( heta) &= \int_{t_0}^{t_1} \|\epsilon^ op 
abla_z f(\mathbf{x}(t),t, heta)\|_2^2 \,\mathrm{d}t \end{aligned}$$

CIFAR10		ImageNet64	
BPD	Time	BPD	Time
3.38	31.84	3.83	64.1

https://arxiv.org/abs/2002.02798

**Continuous Normalizing Flows** 

STEER

Introduces temporal regularization:

$$egin{aligned} \mathbf{x}(t_1) &= \mathbf{x}(t_0) + \int_{t_0}^T f_ heta(\mathbf{x}(t),t) \mathrm{d}t \ &= \mathrm{ODESolve}(\mathbf{x}(t_0),f_ heta,t_0,T) \ &T \sim \mathrm{Uniform}(t_1-b,t_1+b) \ &b < t_1-t_0 \end{aligned}$$

CIFAR10		ImageNet64	
BPD	Time	BPD	Time
3.39	22.24	-	-

https://arxiv.org/abs/2006.10711

### Vikram Voleti

# **3. Image Generation**





https://arxiv.org/abs/1810.01367

### **Continuous Normalizing Flows**

https://arxiv.org/abs/2006.10711

#### Vikram Voleti



# Thank you!

voletiv.github.io