### Latent Stochastic Differential Equations







[Hegde et al., 2018]

Xuechen Li, Leonard Wong, Ricky Chen, Yulia Rubanova, David Duvenaud University of Toronto, Vector Institute





### Irregularly-timed datasets arise all the time



- Most patient data, gene assays irregularly sampled through time.
- All sorts of observation models (likelihoods), not just Gaussian
- Most large parametric models in ML are discrete time: RNNs, HMM, DMM







- **Imputation**: Need to solve original problem, messes with uncertainty

**Binning**: End up averaging many entries per bin, or leaving bins empty



### Latent variable models

- Hidden Markov Models, Deep Markov Models
  - specify p(z), p(x | z) $p(x) = \int p(x | z)p(z)dz$
- Can integrate out z however you want!
  - Variational inference, MCMC lacksquare
  - A neural net or whatever can specify proposal or approx. posterior



https://pyro.ai/examples/dmm.html

• [Krishnan, Shalit, Sontag]



### Latent variable models

- Can use a neural net to guess optimal variational params from data
- Structure of recognition net an implementation detail
  - Only there to speed things up.
  - Just needs to output a normalized distribution over z



https://pyro.ai/examples/dmm.html

# What about continuous time?



- Vector-valued **z** changes in time

• Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t), t)$ 



- Vector-valued **z** changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t), t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





- Vector-valued **z** changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t), t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued **z** changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued **z** changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





t

- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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





- Vector-valued z changes in time
- Time-derivative:  $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t),t)$
- Initial-value problem: given  $\mathbf{z}(t_0)$ , find:

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



## Autoregressive continuous-time?

- ODE-RNN:
  - Between datapoints, dh/dt = f(h(t))
  - At observations,  $h(t)' = g(x_t, h(t))$
- h represents belief state (like in an RNN)
- Separates belief update due to time passing vs seeing data. Good!





### ODE latent-variable model ODE Solve $(z_{t_0}, f, \theta_f, t_0, ..., t_N)$ $\langle z_{t_1} \rangle$ $Z_{t_N}$ $\langle z_{t_i} \rangle$ $z_{t_0}$ $\hat{x}_{t_i}$ $(\hat{x}_{t_1})$ $(\hat{x}_{t_0})$

- z(t) represents true state of system at time t
- Need to approximate posterior p(z\_t0 | x\_t1...)
- Well-defined state at all times, dynamics separate from inference

 $\mathbf{z}_{t_0} \sim p(\mathbf{z}_{t_0})$  $\mathbf{z}_{t_1}, \mathbf{z}_{t_2}, \ldots, \mathbf{z}_{t_N} = \text{ODESolve}(\mathbf{z}_{t_0}, f, \theta_f, t_0, \ldots, t_N)$ each  $\mathbf{x}_{t_i} \sim p(\mathbf{x} | \mathbf{z}_{t_i}, \theta_{\mathbf{x}})$ 





## An ODE latent-variable model

• Can do VAE-style inference with an RNN encoder







Latent trajectories z(t) (latent space)









Slice of vector field (latent space)

Samples from prior (data space)







Latent trajectories z(t) (latent space)









Slice of vector field (latent space)

Samples from prior (data space)



## Mujoco: State versus Belief states

- States are more interpretable than belief states
- True dynamics are deterministic

Truth

Latent ODE

||f(z)||(ODE)

 $||\Delta h||$  (RNN)







Time





Table 4: Test MSE (mean  $\pm$  std) on PhysioNet. Autoregressive models.

Model	Interp ( $\times 10^{-3}$ )	Encouer-accouer models.		
		Model	Interp ( $\times 10^{-3}$ )	Extrap (×10
RNN $\Delta_t$ RNN-Impute RNN-Decay	$3.520 \pm 0.276$ $3.243 \pm 0.275$ $3.215 \pm 0.276$	RNN-VAE Latent ODE (RNN enc.)	$5.930 \pm 0.249$ $3.907 \pm 0.252$	$3.055 \pm 0.1$ $3.162 \pm 0.0$
RNN GRU-D	$3.384 \pm 0.274$	Latent ODE (ODE enc) Latent ODE + Poisson	$2.118 \pm 0.271$ 2.789 ± 0.771	$2.231 \pm 0.0$ $2.208 \pm 0.0$
ODE-KININ (Ours)	$2.301 \pm 0.000$			

Table 5: Test MSE (mean  $\pm$  std) on PhysioNet. Fncoder\_decoder models



### Poisson Process Likelihoods

 $\log p(t_1,\ldots,t_N|t_{\text{start}},t_{\text{end}})$ 

$$=\sum_{i=1}^N \log \lambda(\mathbf{z}(t_i)) - \int_{t_{ ext{start}}}^{t_{ ext{end}}} \lambda(\mathbf{z}(t)dt)$$

- Model joint p(obs, time) instead of p(obs | time)
- Non-intervention model









- Latent ODEs for Irregularly-Sampled Time Series
  - Yulia Rubanova, Ricky T. Q. Chen, David Duvenaud
- <u>https://github.com/</u> YuliaRubanova/latent ode







### Limitations of Latent ODEs ODE Solve $(z_{t_0}, f, \theta_f, t_0, ..., t_N)$ $z_{t_{1}}$ $\langle z_{tN} \rangle$ $(z_{t_0})$ $\langle z_{t_i} \rangle$ sequence length $\hat{x}_{t_1}$ $(\hat{x}_{tN})$ $(\hat{x}_{t_i})$ $(\hat{x}_t)$

- Deterministic dynamics!
  - State size grows with
- Special time t0, only reason about z\_t0



### Let's be like a Deep Markov Model

Nonlinear latent variable with noise at each step:

$$z_{t+1} = z_t + f_{\theta}(z_t) + \epsilon$$

- Could add more steps between observations.
- Infinitesimal limit some sort of stochastic ODE...?



https://pyro.ai/examples/dmm.html

## Stochastic Differential Equations

 $\frac{dz}{dt} = f(z(t)) + \text{``}\epsilon\text{''}$ 

### $dz = f(z(t))dt + \sigma(z(t))dB(t)$

Implicit distribution over functions



- natural fit for many small, unobserved interactions:
  - motion of molecules in a liquid
  - allele frequencies in a gene pool
  - prices in a market
- Interactions don't need to be Gaussian; as long as CLT kicks in, you get Brownian motion
- Let's put neural networks into SDE dynamics and fit giant SDE models to everything!

### Life is an SDE

 $dz = f_{\theta}(z(t))dt + \sigma_{\theta}(z(t))dB(t)$ 

### Neural Ordinary Differential Equations

Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud

### Neural Stochastic Differential Equations: Deep Latent Gaussian Models in the Diffusion Limit

Belinda Tzen, Maxim Raginsky

### Neural Stochastic Differential Equations

Stefano Peluchetti, Stefano Favaro

### Neural Jump Stochastic Differential Equations

Junteng Jia, Austin R. Benson



incile een



 Tzen + Raginski: Deep LVMs become SDEs in the limit. Variational inf framework. Forwardmode autodiff.

### Neural Ordinary **Differential Equations**

Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud

### Neural Stochastic Differential Equations: Deep Latent Gaussian Models in the Diffusion Limit

Belinda Tzen, Maxim Raginsky

### **Neural Stochastic Differential Equations**

Stefano Peluchetti, Stefano Favaro

### **Neural Jump Stochastic Differential Equations**

Junteng Jia, Austin R. Benson



### Intellocent



- Tzen + Raginski: Deep LVMs become SDEs in the limit. Variational inf framework. Forwardmode autodiff.
- Peluchetti + Favaro: Worked out SDE corresponding to infinitelydeep convnets with uncertain weights

### Neural Ordinary **Differential Equations**

Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud

### Neural Stochastic Differential Equations: Deep Latent Gaussian Models in the Diffusion Limit

Belinda Tzen, Maxim Raginsky

### **Neural Stochastic Differential Equations**

Stefano Peluchetti, Stefano Favaro

### **Neural Jump Stochastic** Differential Equations

Junteng Jia, Austin R. Benson



Intellocent



- Tzen + Raginski: Deep LVMs become SDEs in the limit. Variational inf framework. Forwardmode autodiff.
- Peluchetti + Favaro: Worked out SDE corresponding to infinitelydeep convnets with uncertain weights
- Jia + Benson: Added countably many discrete jumps to latent ODEs

### Neural Ordinary **Differential Equations**

Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud

### Equations: Deep Latent Gaussian

### **Differential Equations**

Stefano Peluchetti, Stefano Favaro

### **Neural Jump Stochastic** Differential Equations

Junteng Jia, Austin R. Benson

Neural Stochastic Differential Models in the Diffusion Limit Belinda Tzen, Maxim Raginsky Neural Stochastic






• Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box* variational inference for stochastic differential equations.

- Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box* variational inference for stochastic differential equations.
- Pashupati Hegde, Markus Heinonen, Harri Lähdesmäki, and Samuel Kaski. Deep learning with differential gaussian process flows.

- Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box variational inference for stochastic differential equations.*
- Pashupati Hegde, Markus Heinonen, Harri Lähdesmäki, and Samuel Kaski. *Deep learning with differential gaussian process flows.*
- Markus Heinonen, Cagatay Yildiz, Henrik Mannerström, Jukka Intosalmi, and Harri Lähdesmäki. *Learning unknown ODE models with gaussian processes.*

- Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box variational inference for stochastic differential equations.*
- Pashupati Hegde, Markus Heinonen, Harri Lähdesmäki, and Samuel Kaski. *Deep learning with differential gaussian process flows.*
- Markus Heinonen, Cagatay Yildiz, Henrik Mannerström, Jukka Intosalmi, and Harri Lähdesmäki. *Learning unknown ODE models with gaussian processes.*
- C. Garcıa, A. Otero, P. Felix, J. Presedo, and D. Marquez. Nonparametric estimation of stochastic differential equations with sparse Gaussian processes.

- Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box* variational inference for stochastic differential equations.
- Pashupati Hegde, Markus Heinonen, Harri Lähdesmäki, and Samuel Kaski. Deep learning with differential gaussian process flows.
- Markus Heinonen, Cagatay Yildiz, Henrik Mannerström, Jukka Intosalmi, and Harri Lähdesmäki. Learning unknown ODE models with gaussian processes.
- C. Garcia, A. Otero, P. Felix, J. Presedo, and D. Marguez. Nonparametric estimation of stochastic differential equations with sparse Gaussian processes.

 All use Euler discretizations. Not clear what limiting algorithm is (e.g. enforces invariants?), and not memory-efficient.

- Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box* variational inference for stochastic differential equations.
- Pashupati Hegde, Markus Heinonen, Harri Lähdesmäki, and Samuel Kaski. Deep learning with differential gaussian process flows.
- Markus Heinonen, Cagatay Yildiz, Henrik Mannerström, Jukka Intosalmi, and Harri Lähdesmäki. Learning unknown ODE models with gaussian processes.
- C. Garcia, A. Otero, P. Felix, J. Presedo, and D. Marguez. Nonparametric estimation of stochastic differential equations with sparse Gaussian processes.

- All use Euler discretizations. Not clear what limiting algorithm is (e.g. enforces invariants?), and not memory-efficient.
- Not even going to discuss methods that require solving a PDE - not scalable.



- Thomas Ryder, Andrew Golightly, A Stephen Mc-Gough, and Dennis Prangle. *Black-box* variational inference for stochastic differential equations.
- Pashupati Hegde, Markus Heinonen, Harri Lähdesmäki, and Samuel Kaski. Deep learning with differential gaussian process flows.
- Markus Heinonen, Cagatay Yildiz, Henrik Mannerström, Jukka Intosalmi, and Harri Lähdesmäki. Learning unknown ODE models with gaussian processes.
- C. Garcia, A. Otero, P. Felix, J. Presedo, and D. Marguez. Nonparametric estimation of stochastic differential equations with sparse Gaussian processes.

- All use Euler discretizations. Not clear what limiting algorithm is (e.g. enforces invariants?), and not memory-efficient.
- Not even going to discuss methods that require solving a PDE - not scalable.
- We want to use adaptive, (high-order?) SDE solvers.



# How to fit ODE params?

 $\frac{\partial L}{\partial \theta} = ?$ 

 $L(\theta) = L\left(\int_{t_0}^{\tau_1} f(\mathbf{z}(t), t, \theta) dt\right)$ 

# How to fit ODE params? $L(\theta) = L\left(\int_{t_0}^{t_1} f(\mathbf{z}(t), t, \theta) dt\right)$ $\frac{\partial L}{\partial A} = ?$

- Alexey Radul: Approximate the derivative, don't differentiate the approximation!

### Don't backprop through solver: High memory cost, extra numerical error



Standard Backprop:



 $\frac{\partial L}{\partial \theta} = \sum_{t} \frac{\partial L}{\partial \mathbf{z}_{t}} \frac{\partial f(\mathbf{z}_{t}, \theta)}{\partial \theta}$ 

Standard Backprop:



 $\frac{\partial L}{\partial \theta} = \sum_{t} \frac{\partial L}{\partial \mathbf{z}_{t}} \frac{\partial f(\mathbf{z}_{t}, \theta)}{\partial \theta}$ 

Adjoint sensitivities: (Pontryagin et al., 1962):

 $\frac{\partial}{\partial t} \frac{\partial L}{\partial \mathbf{z}(t)} = \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \mathbf{z}}$  $\frac{\partial L}{\partial \theta} = \int_{t_1}^{t_0} \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \theta} dt$ 

Adjoint sensitivities: (Pontryagin et al., 1962):

 $\frac{\partial}{\partial t} \frac{\partial L}{\partial \mathbf{z}(t)} = \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \mathbf{z}}$  $\frac{\partial L}{\partial \theta} = \int_{t_1}^{t_0} \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \theta} dt$ 

• Can build adjoint dynamics with autodiff, compute gradients with second ODE solve:

Adjoint sensitivities: (Pontryagin et al., 1962):

> $\frac{\partial}{\partial t} \frac{\partial L}{\partial \mathbf{z}(t)} = \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \mathbf{z}}$  $\frac{\partial L}{\partial \theta} = \int_{t_1}^{t_0} \frac{\partial L}{\partial \mathbf{Z}(t)} \frac{\partial f(\mathbf{Z}(t), \theta)}{\partial \theta} dt$

 Can build adjoint dynamics with autodiff, compute gradients with second ODE solve:

def f\_aug([z, a, d], t): return [f, -a\*df/dz, -a\*df/dθ)

Adjoint sensitivities: (Pontryagin et al., 1962):

> $\frac{\partial}{\partial t} \frac{\partial L}{\partial \mathbf{z}(t)} = \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \mathbf{z}}$  $\frac{\partial L}{\partial \theta} = \int_{t_1}^{t_0} \frac{\partial L}{\partial \mathbf{Z}(t)} \frac{\partial f(\mathbf{Z}(t), \theta)}{\partial \theta} dt$

- Can build adjoint dynamics with autodiff, compute gradients with second ODE solve:
- def f\_aug([z, a, d], t): return [f, -a\*df/dz, -a\*df/dθ)
- $[z0, dL/dz(t0), dL/d\theta] =$ ODESolve(f\_aug, [z(t1), dL/dz(t1), 0], t1, t0)

Adjoint sensitivities: (Pontryagin et al., 1962):

> $\frac{\partial}{\partial t} \frac{\partial L}{\partial \mathbf{z}(t)} = \frac{\partial L}{\partial \mathbf{z}(t)} \frac{\partial f(\mathbf{z}(t), \theta)}{\partial \mathbf{z}}$  $\frac{\partial L}{\partial \theta} = \int_{t_1}^{t_0} \frac{\partial L}{\partial \mathbf{Z}(t)} \frac{\partial f(\mathbf{Z}(t), \theta)}{\partial \theta} dt$

- No need to store activations, just run dynamics backwards from output.
- Easy to run ODE backwards, just run negate dynamics and time:
  - back\_f(z, t) = -f(z, -t)





### Why not repeat same trick?

# not apply same adjoint method?

If an SDE is just "an ODE with noise", why

- Reparameterization trick: Use same noise from forward pass on reverse pass
- Infinite reparameterization trick: Use same Brownian motion sample on forward and reverse passes.

### Need to store noise





- Can 'zoom in' arbitrarily close at any point

t

• splittable random seed ensures all entire sample is consistent

### type UnitInterval = Real

```
bmIter :: (Key, Real, Real, UnitInterval) -> (Key, Real, Real, UnitInterval)
bmIter (key, y, sigma, t) =
  (kDraw, kL, kR) = splitKey3 key
 t' = abs (t - 0.5)
  y' = sigma * randn kDraw * (0.5 - t')
 key' = select (t > 0.5) kL kR
  (key', y + y', sigma / sqrt 2.0, t' * 2.0)
sampleBM :: Key -> UnitInterval -> Real
sampleBM key t =
  (_, y, _, _) = fold (key, 0.0, 1.0, t) for i::10. bmIter
```

```
У
```





• Me: Let's just slap negative signs on everything and hope for the best



- Me: Let's just slap negative signs on everything and hope for the best
- Xuechen and Leonard: What does that even mean? Much later: Actually we proved that's correct.



- Me: Let's just slap negative signs on everything and hope for the best
- Xuechen and Leonard: What does that even mean? Much later: Actually we proved that's correct.
- Builds on results from Kunita 2019.



- Me: Let's just slap negative signs on everything and hope for the best
- Xuechen and Leonard: What does that even mean? Much later: Actually we proved that's correct.
- Builds on results from Kunita 2019.
- Adjoint formula is analogous to ODE.



### Generalize adjoint to diffusion func

- Dynamics already known
- Diffusion adjoint almost the same as drift adjoint.
- Just more vector-Jacobian products!

def drift\_adjoint(augmented\_state, t, flat\_args): y, y\_adj, args\_adj = unpack(augmented\_state) fval, vjp\_all = vjp(flat\_f, y, t, flat\_args) f\_vjp\_a, f\_vjp\_t, f\_vjp\_args = vjp\_all(-y\_adj) return np.concatenate([fval, f\_vjp\_a, f\_vjp\_args])

def diffusion\_adjoint\_prod(augmented\_state, t, args, v): y, y\_adj, arg\_adj = unpack(augmented\_state) gval, vjp\_g = vjp(flat\_g, y, t, args) vjp\_a, vjp\_t, vjp\_args = vjp\_g(-y\_adj \* v) **return** np.concatenate([gval  $* v_{i} v_{j}p_{a_{i}} v_{j}p_{a}]$ )

$$\begin{split} \tilde{A}_{s,t}(z) = \nabla \mathcal{L}(z) + \int_{s}^{t} \nabla b(\check{\Psi}_{r,t}(z), r)^{\top} \tilde{A}_{r,t}(z) \, \mathrm{d}r + \\ \sum_{i=1}^{m} \int_{s}^{t} \nabla \sigma_{i}(\check{\Psi}_{r,T}(r), u)^{\top} \tilde{A}_{r,t}(z) \circ \check{\mathrm{d}}W_{r}^{i}, \end{split}$$







### Time complexity (fixed-step)

Method Time Memory  $\mathcal{O}(LD)$ Forward pathwise [14, 60] $\mathcal{O}(1)$ Backprop through solver [11]  $\mathcal{O}(L)$  $\mathcal{O}(L)$  $\mathcal{O}(L \log L)$ Stochastic adjoint (ours)  $\mathcal{O}(1)$ 

- Just solving forwards costs O(L) time.
- Okay! Now we can fit SDE paths to data...

Time more like O(L) when dynamics are expensive

## Need Latent (Bayesian) SDE

- Can't just fit SDE to maximize likelihood optimal solution has no randomness and just hugs data
- Define prior and approx. posterior implicitly:

### $dz_p = f_{\theta}(z(t))dt + \sigma_{\theta}(z(t))dB(t)$







### Variational inference

• Define prior and approx. posterior implicitly:

 $dz_p = f_{\theta}(z(t))dt + \sigma_{\theta}(z(t))dB(t)$  $dz_{q} = f_{\phi}(z(t))dt + \sigma_{\theta}(z(t))dB'(t)$  $u(t) = \left| \frac{f_{\theta}(z(t)) - f_{\phi}(z(t))}{\sigma_{\theta}(z(t))} \right|_{2}^{2}$  $\mathcal{L}_{\mathrm{VI}} = \mathbb{E}\left[\frac{1}{2}\int_0^T |u(Z_t, t)|^2 \, \mathrm{d}t - \sum \log p(y_{t_i}|z_{t_i})\right] \quad \text{(a) Inference}$ i=1



### Latent SDEs: An unexplored model class

- Define implicit prior over functions with neural nets for dynamics (like GP hyperparams)
- Define implicit approximate posterior through neural nets for dynamics (variational params)
- Define observation likelihoods. Anything differentiable wrt latent state will work (e.g. text models!)
- Train everything jointly with ADVI. Should scale to millions of params. Can use Milstein solver for diagonal noise.



samples from learned prior dynamics





### Early latent SDE results 1.5 -

0.5 -

1.0 -

- Ľ 0.0 -• OU prior, Laplace likelihood -0.5 -
  - -1.0 -
  - -1.5 -





0.00 0.25 0.50 0.75 1.00 1.25 1.50 1.75 2.00

### Early latent SDE results 1.5 -

0.5 -

1.0 -

- Ľ 0.0 -• OU prior, Laplace likelihood -0.5 -
  - -1.0 -
  - -1.5 -





0.00 0.25 0.50 0.75 1.00 1.25 1.50 1.75 2.00

# Early latent SDE results

- OU prior,
   Laplace likelihood
- Inference problem is more <sup>>\*</sup> 0.0 local in time than for ODE (recognition net can steer posterior sample towards data)
   -0.5 --1.0 -

-1.5 -

1.5 -

1.0 -

0.00 0.25 0.50 0.75 1.00 1.25 1.50 1.75 2.00 t



# Early latent SDE results

- OU prior,
   Laplace likelihood
- Inference problem is more <sup>>\*</sup> 0.0 local in time than for ODE (recognition net can steer posterior sample towards data)
   -0.5 --1.0 -

-1.5 -

1.5 -

1.0 -

0.00 0.25 0.50 0.75 1.00 1.25 1.50 1.75 2.00 t



### Early latent SDE results: Mocap

 50D data, 6D latent space, sharing dynamics and recognition params across time series (11000 params)



Method	Test MSE
npODE [18]	$22.96^{+}$
NeuralODE $[4]$	$22.49 \pm 0.$
$ODE^2VAE$ [61]	$10.06 \pm 1.$
$ODE^2VAE-KL$ [61]	$8.09 \pm 1.9$
Latent ODE $[4, 50]$	$5.98\pm0.2$
Latent SDE (this work)	$4.03 \pm 0.2$


## Early latent SDE results: Mocap

 50D data, 6D latent space, sharing dynamics and recognition params across time series (11000 params)



Method	Test MSE
npODE [18]	$22.96^{+}$
NeuralODE $[4]$	$22.49 \pm 0.$
$ODE^2VAE$ [61]	$10.06 \pm 1.$
$ODE^2VAE-KL$ [61]	$8.09 \pm 1.9$
Latent ODE $[4, 50]$	$5.98\pm0.2$
Latent SDE (this work)	$4.03 \pm 0.2$



## Limitations

- SDE solvers much lower-order convergence than ODE solvers
  - (e.g. Milstein order 1 vs RK4)
- Non-diagonal noise requires Levy areas
  - Diagonal noise requires funny parameterization
- Need jump-style noise? (e.g. hit by a car)
- Not scalable in input dimension (diffusions)?

# SDES VS GPS

- Distinct sets of priors over functions
- Easy to construct non-Gaussian SDE



Line means "can be used to construct", but not "contains"

### Future work:

- Skipping short-scale dynamics between observations (mixes back to prior)
- Infinitely deep Bayesian neural networks
- Code in a week or two (prototype)





#### Xuechen Li, Leonard Wong, Ricky Chen, Yulia Rubanova, David Duvenaud







### Thanks!



