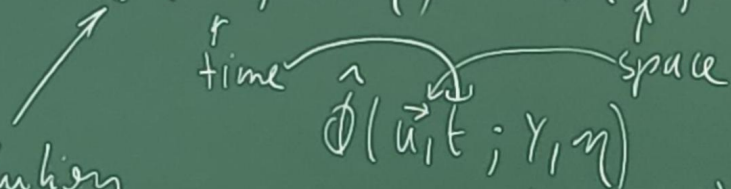




- PINNs address the case where  $X$  is a function surrogate

$$\hat{\Phi}(t; \gamma, \eta) \quad \hat{\Phi}(\vec{u}; \gamma, \eta)$$



$\Phi(t)$  is solution

of a ordinary differential eq

$$\hat{\Phi}(\vec{u}; t; \gamma, \eta)$$

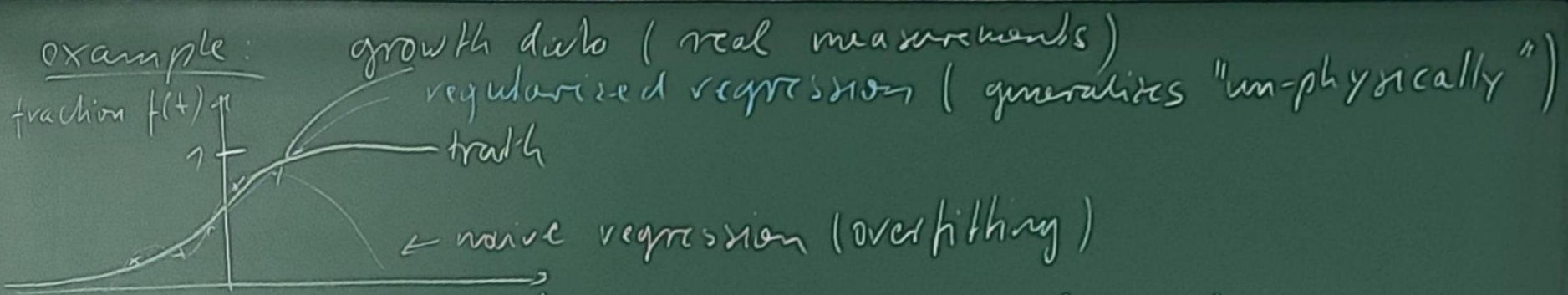
$\Phi(\vec{u}; \gamma, \eta)$  and  $\Phi(\vec{u}; t; \gamma, \eta)$

solutions of partial differential eq

- physical prior knowledge: formula for the ODE or PDE

- why is this useful?

- actual simulation may be too expensive
- requires a lot less training data than traditional SBI
- find a model that fits experimental data



idea: use physical knowledge as a problem-specific regularizer

• here: growth follows "logistic rule".  $\frac{df(t)}{dt} = \lambda f(t)(1-f(t))$   
 $\lambda$  may or may not be known

• special case of SIR eq.

$$\frac{dS}{dt} = -\lambda \frac{S \cdot I}{N} \quad \frac{dI}{dt} = \lambda \frac{S \cdot I}{N} - \mu I \quad \frac{dR}{dt} = \mu I$$

set  $\mu=0$  and reformulate in terms of fractions,  $[S] = \frac{S}{N}$  etc.

$$\frac{d[S]}{dt} = -\lambda [S][I] \quad \frac{d[I]}{dt} = \lambda [S][I] \quad \text{s.t. } [S] + [I] = 1$$



$$f(t) = [I](t) = \sigma(\lambda(t-t_0)) = \frac{1}{1 + \exp(-\lambda(t-t_0))}$$

• PINN approach to learning  $\hat{f}(t; \lambda, t_0, \eta)$

- two kinds of points (later, 3)  $\leftarrow$  observation noise

• data points  $\hat{=}$  TS where  $f(t)$  is (approximately) known  
(at least data has initial condition  $t=0$ )

• collocation points points where we apply the regularizer

$t \in CP$  check if ODE is fulfilled at  $t$

$\rightarrow$  regularization term

$$h_{ODE} = \frac{1}{M} \sum_{m=1}^M \left( \frac{df(t_m)}{dt} - \lambda f(t_m)(1-f(t_m)) \right)^2$$

$\stackrel{!}{=} 0$  if ODE is fulfilled

(ideally  $M \rightarrow \infty$   
but finite  $M$  is  
sufficient in practice)

data term: 
$$h_{data} = \frac{1}{N} \sum_{i=1}^N (f_i - f(t_i))^2$$





general case for ODEs:  $x(t), \dot{x}(t) = \frac{dx(t)}{dt}, \ddot{x}(t) = \frac{d\dot{x}(t)}{dt} = \frac{d^2x(t)}{dt^2}$

- ODE  $\text{ODE}(x, \dot{x}, \ddot{x}; t, Y) = 0$

① define TS =  $\{(t_i, x_i, \dot{x}_i, \ddot{x}_i)\}_{i=1}^N$  and allocation set  $\{CS = \{t_m \sim p(t_m)\}_{m=1}^M$   
both real and simulated optional [or re sample CS for every batch]  
if Y is known

① from  $\tau = 1, \dots, \tau_{max}$

① do a gradient step of  $\mathcal{L}_{data} \sim \mathcal{L}_{data} + \mathcal{L}_{ODE} \sim \mathcal{L}_{ODE} + \mathcal{L}_{bound} \sim \mathcal{L}_{bound}$   
(usually ADAM)  
optional include a loss for boundary conditions (if applicable)

ex logistic ODE  $0 \leq f(t) \leq 1$

$$\mathcal{L}_{bound} = \frac{1}{B} \sum_{b=1}^B \left( \text{ReLU}(-f(t_b)) + \text{ReLU}(f(t_b) - 1) \right)^2$$

general form of loss terms

$$L_{data} = \frac{1}{N} \sum_{i=1}^N (f_i - f(t_i))^2$$

$$L_{ODE} = \frac{1}{M} \sum_{m=1}^M \text{ODE}(f(t_i), t_i, t_m, Y)^2$$

$$L_{bound} = \frac{1}{B} \sum_{b=1}^B B(f(t_i), t_i, t_{b1})^2 \quad \text{boundary constraint functions}$$

benefits of using neural networks

- universal approximators (if big enough) and good convergence in practice  
 $\Rightarrow$  can in principle learn any  $f^*(t, Y)$
- derivatives  $f'(t)$ ,  $f''(t)$  easily computable by autodiff

disadvantage:

- for each set of data points (incl. initial conditions) and parameters  $X$  learning must start from scratch  $\Rightarrow$  no amortization / generalization  
problem is currently addressed  $\Rightarrow$  later

tricks to improve performance:

- use  $\tanh(u)$  activation or recently

$$\text{GReLU}(u) = \frac{u}{2} \cdot \sqrt{1 + \text{erf}\left(\frac{u}{\sqrt{2}}\right)}$$

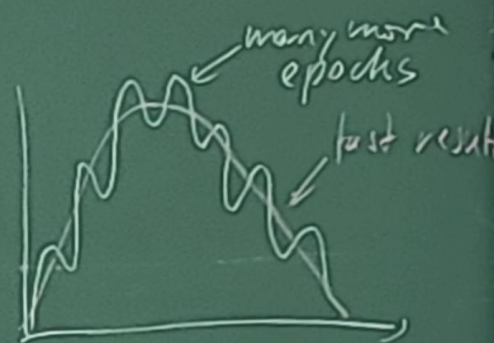
(CDF of standard normal)



- standardize the problem (similar to scaling features to unit variance in standard regression)  
transform ODE into a "dimensionless" form

- use fraction  $[I] = \frac{I}{N}$  instead of counts  $I \Rightarrow 0 \leq f(t) \leq 1$

- choose units of free parameters  $\gamma$  "cleverly"  
[all  $\gamma_j$  should be  $O(1)$  ?]

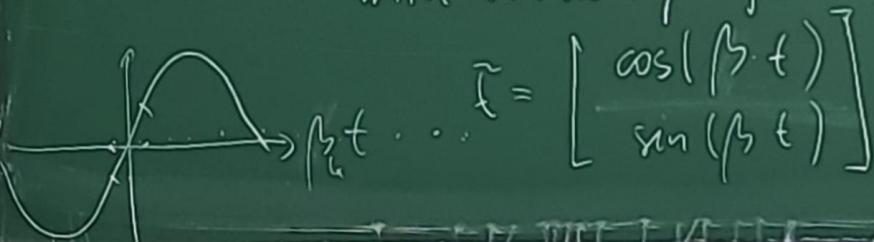


- network architecture

- fully connected networks
- random Fourier features

• was shown [Nahaman et al. 2016] that low frequency behavior of  $f(t)$  converges much faster than high frequency behavior

$\Rightarrow$  convert features (here: time coordinate) to the Fourier domain with random projections



$$\tilde{x} = \begin{bmatrix} \cos(\beta_k t) \\ \sin(\beta_k t) \end{bmatrix}$$

$$\beta = [\beta_1, \beta_k] \quad \beta_k \sim N(0, \sigma^2)$$

$$\sigma^2 \in [1, 100]$$



