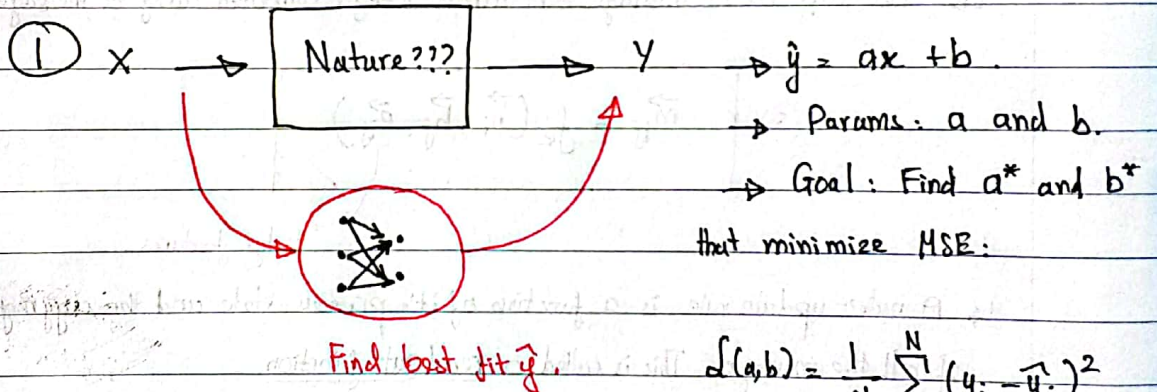Personal Note #8: Graph Neural Diffusion.
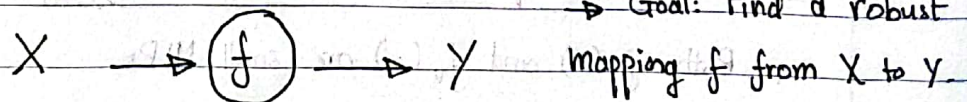
① Prerequisites: ODE and Neural ODE.
ⓐ Understanding Ordinary Differential Equation.

→ Traditional data modelling method (regression): Given a bunch of paired input and output $\{(x_i, y_i)\}_{i=1}^{N}$, find a function $f: X \to Y$ to approx output directly.

→ The problem with this traditional paradigm is that the data is generated by nature factors (economic, social, ...). We have no reliable way to remodel the sophisticated natural generation process. So, we treat nature like a black box & by pass it.



① X → [Nature???] → Y     → $\hat{y} = ax + b$.
                                    → Params: a and b.
                                    → Goal: Find $a^*$ and $b^*$
Find best fit $\hat{y}$.          that minimize MSE:

$$\mathcal{L}(a,b) = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

② Using ODE:
                                    → Goal: Find a robust
X → ⓕ → Y                        mapping f from X to Y.

ⓑ Solve an ODE numerically: Euler's Method.
ⓧ Modelling rate of change.
  ↳ So, instead of modelling the              function itself directly,
  we model the rate of change.
  ↳ For every continuous differentiable function f, its rate of change
  is defined as the derivative:

$$f' = \frac{df}{dt} \text{ (or } dx, \text{ depending on input).}$$

## ⊛ ODE basic form.

↳ In the basic form of an ODE, we allow the derivative to depend not only on $x$ but also on $y$. This allows modelling flexibility.

↳ When integrating the derivative, the integral depends on at least 1 constant variable. To solve for this constant, we need an initial condition $(x_0, y_0)$ (initial value problem).

⊛
$$y'(x) = f(x,y) \text{ where } y(x_0) = y_0.$$
⟶ ⁺ Typical IVP.

## ⊛ Euler's Method.

↳ Given a derivative of a function and an initial value :

$$\begin{cases} \dfrac{dy}{dx} = f(x,y) \\[2mm] y(x_0) = y_0. \end{cases}$$

↳ We know that $\dfrac{dy}{dx}(x_0, y_0)$ defines the slope of the line tangent to the curve at $(x_0, y_0)$.
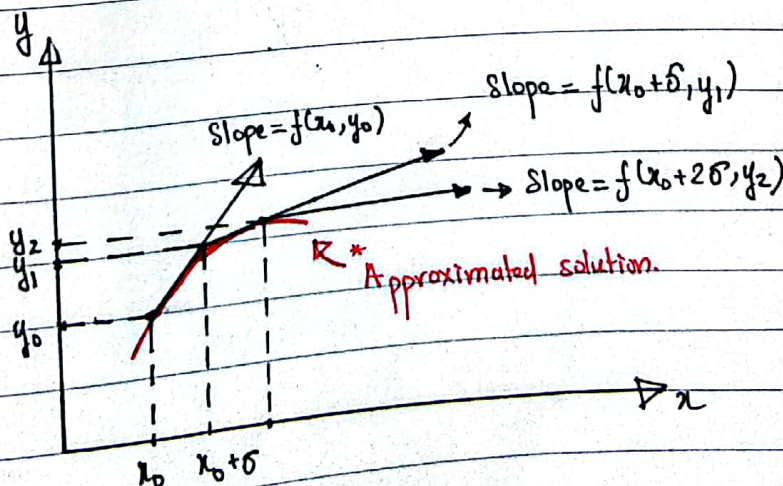
↳ Euler's method is derived from the basic form of tangential approximation :

⊛
$$y(x + \delta) \approx y(x) + \delta \frac{dy}{dx}$$

⟶ ⁺ $\delta$ is a fixed step size. The smaller the step, the more accurate the approximation.



Slope $= f(x_0 + \delta, y_1)$

Slope $= f(x_1, y_0)$

Slope $= f(x_0 + 2\delta, y_2)$

R * Approximated solution.

$x_0 \quad x_0 + \delta$

⌐→ In simple terms, Euler's method starts at an initial condition $(x_0, y_0)$ then recursively increase $x_0$ by a fixed step size $\delta$ and evaluate the slope at each step to derive the approximated solution (curve).

(*)
$$\begin{cases} x_{n+1} = x_n + \delta \\ y_{n+1} = y_n + \delta f(x_n, y_n) \end{cases}$$

(*) Optimising the ODE.

⌐→ Suppose we want to model a linear regression model $\hat{y}(x) = ax + b$.
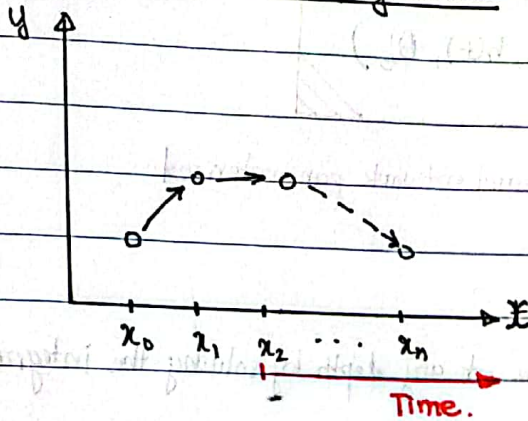
$$\to \frac{dy}{dx} = a.$$

→ We have a set of input X and output Y with N data points.

⌐→ To model $\hat{y}$ using an ODE,

1. Choose an initial value $(x_0, y_0)$ and a fixed stepsize $\delta$.
2. Choose k data points to evaluate at $(k < N)$.
3. Add additional point from X so that Euler's method can be computed at a regular interval $\delta$. This is my computation trajectory.
4. Evaluate the trajectory at the chosen k points.
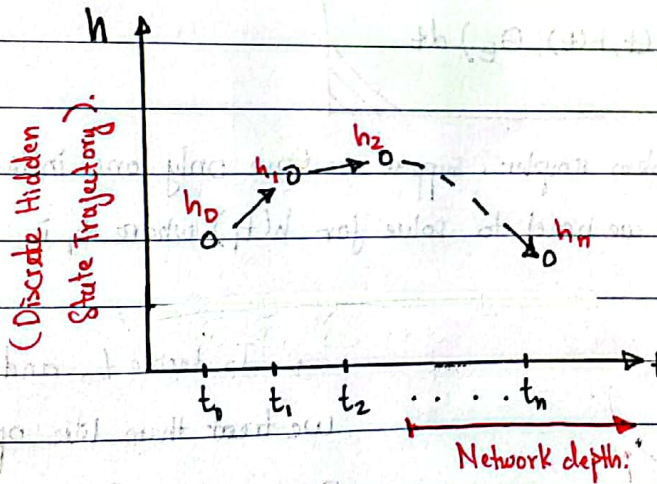5. Compute loss & gradient descent.

# © Neural ODE.

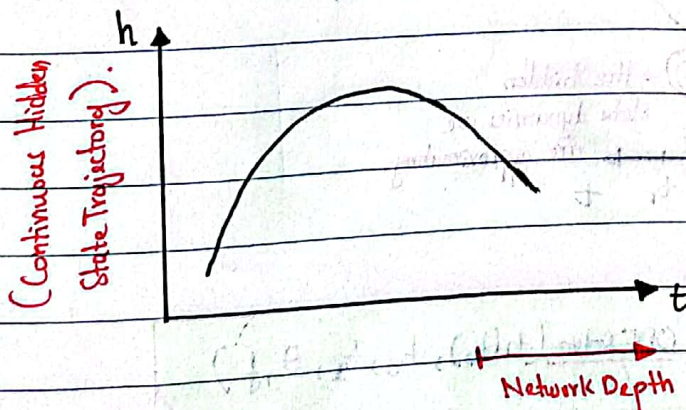## ⊛ Continuous Hidden State Dynamics.



→ In Euler's method, we define the computation trajectory from some initial point $(x_0, y_0)$ and recursively evaluating the ODE at fixed step size.

→ This defines a dynamic of $y$ with respect to time.

(Discrete Hidden State Trajectory).



→ The key idea of neural network is that as we add layers to our neural network, we add hidden states between input & output.

→ We can view the transformation made by each layer a dynamics of hidden state through network depth.

(Continuous Hidden State Trajectory).



→ When we take the continuous limit of hidden state w.r.t network depth, we smooth out computation trajectory.

→ The hidden state can be evaluated at ANY depth. (Main idea of neural ODE)

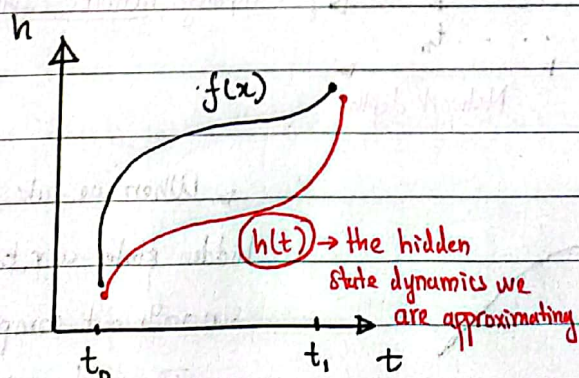↳ In Neural ODE, we parameterise this hidden state dynamics by :

$$ (\ast) \quad \frac{d\,h(t)}{dt} = f(t, h(t), \theta_t) $$

↳ $f(t, h(t), \theta_t)$ is a neural network parameterized by $\theta_t$ at layer $t$.

↳ We can evaluate the hidden state at any depth by solving the integral

$$ (\ast) \quad h(t) = \int f(t, h(t), \theta_t)\, dt $$

numerically. To make the problem simpler, suppose we have only one initial observation $(t_0, h(t_0))$ and we need to solve for $h(t_1)$ where $t_1$ is at the end of the trajectory.
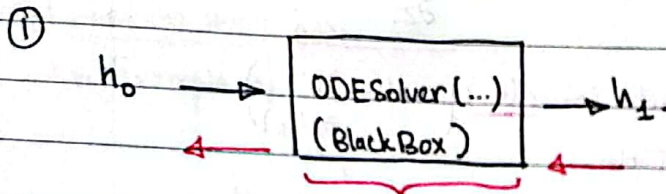
→ To decide $t_0$ and $t_1$ we treat them like optimizable params like $\theta$.



$f(x)$

$h(t)$ → the hidden state dynamics we are approximating.

$$ (\ast) \quad \hat{y} = h(t_1) = \underline{ODE\,solve}\,(h(t_0), t_0, t_1, \theta, f) $$

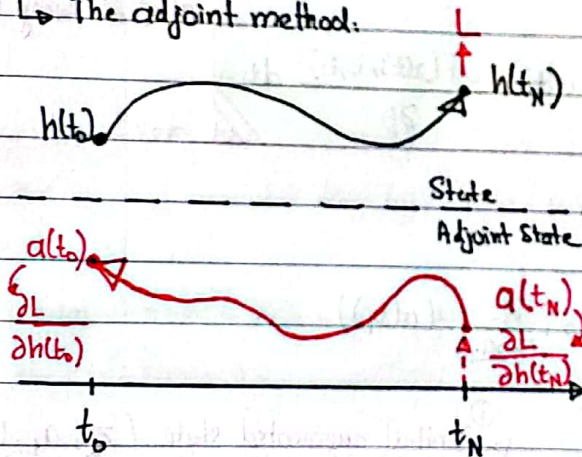*A simple numerical ODE solver. (Euler, Hyper Euler, ...).

## (*) Backpropagating through depth.

↳ The output of the ode net is computed by a "black box" ODE Solver to go from one state $(h_0)$ to another $(h_1)$. Then how do we backprop thru this "black box"?

① 
$$h_0 \longrightarrow \boxed{\begin{array}{c} \text{ODESolver }(\dots) \\ (\text{Black Box}) \end{array}} \longrightarrow h_1.$$

* How do we backward pass
through this black box?

↳ The adjoint method:

→ The difficult thing about training a continuous depth network is back-propagating through the black-box ODE solver.

→ We do this using the adjoint sensitivity method.

$$L$$
$$h(t_N)$$
$$h(t_0)$$
State
Adjoint State
$$a(t_0)$$
$$\frac{\partial L}{\partial h(t_0)}$$
$$a(t_N),$$
$$\frac{\partial L}{\partial h(t_N)}$$
$$t_0 \qquad\qquad t_N$$

↳ Consider optimizing a scalar-valued function $L(\cdot)$:

(*) $$L(z(t_1)) = L\left[ \text{ODESolve}\left( z(t_0), f, t_0, t_1, \theta \right) \right]$$

→ * $z(t)$ = hidden state through time, we define it as $h(t)$ above, just a diff notation.

→ To optimize $L$, The first step is to determine $\frac{\partial L}{\partial z(t)}$. This is called the

adjoint state: (*) $$a(t) = \frac{\partial L}{\partial z(t)}.$$

KOKUYO

↳ The dynamics of the adjoint state is defined by another ODE:

$$\frac{da(t)}{dt} = -a(t)^T \frac{\partial f(z(t),t,\theta)}{\partial z}$$

&* we can solve for $a(t_0)$ (which is $\frac{\partial L}{\partial z(t_0)}$)
by calling ODESolve.

↳ Computing $\frac{\partial L}{\partial \theta}$ requires a third integral which depends on $a(t)$ and $z(t)$.

↳* This is why we need
$a(t)$ at every backward step.

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

↳ Back prop: pseudocode:

Initialization: $\theta, t_0, t_1, z(t_1), \frac{\partial L}{\partial z(t_1)} (a(t_1))$

$s_0 = \left[z(t_1), \frac{\partial L}{\partial z(t_1)}, 0_{|\theta|}\right]$  → ① Initial augmented state $(z_1, a_1, 0)^T$

def aug_dynamics $([z(t), a(t), \cdot], t, \theta)$:  → ② Dynamics of aug state.
    return $\left[f(z(t), t, \theta), -a(t)^T \frac{\partial f}{\partial z}, -a(t)^T \frac{\partial f}{\partial \theta}\right]$  $\left(\frac{dz}{dt}, \frac{da}{dt}, \frac{dL}{d\theta dt}\right)^T$

$\left[z(t_0), a(t_0), \frac{\partial L}{\partial \theta}\right]$ = ODESolve $(s_0, \text{aug\_dynamics}, t_1, t_0, \theta)$

return $a(t_0), \frac{\partial L}{\partial \theta}$ #