Consider solving

$$\frac{\partial u}{\partial t} = F(u, Du, x, t)$$

$$u(x,0) = u_0(x)$$

ex: $u_0 = u_0(x) \in L^1(0,1)$

$$u(0) = u(1) = 0$$

guess

$$u(x,t) = \sum_{n=0}^{\infty} c_n \sin n\pi x$$

truncate at $N$ & minimize

$$\min_{\vec{c} \in \mathbb{R}^N} \|u - u_N\|_2^2$$

Solution:

$$C_n = 2 \langle u_0, \sin n\pi x \rangle_{L^2(0,1)}$$

This was a baller parametrization!

Here's a different approach...
Use a neural network!
......

Let $\quad r_\theta = \dfrac{\partial u_\theta}{\partial t} - \dfrac{\partial^2 u_\theta}{\partial x^2}$

$$L_1 \quad = \quad \frac{1}{N_1} \sum_{i=1}^{N_1} |r_\theta(x_i, t_i)|^2$$

$$L_2 \quad = \quad \frac{1}{N_2} \sum_{i=1}^{N_2} |u_\theta(x_i, 0) - u_0(x_i)|^2$$

$$L_3 \quad = \quad \frac{1}{N_3} \sum_{i=1}^{N_3} |u_\theta(0, t_i)|^2 + |u_\theta(1, t_i)|^2$$

$$M = \underline{\sum L_i}$$

$$\min_{\theta \in \mathbb{R}^M} \sum_{i=1}^{3} \alpha_i L_i$$

How we learn!

(*) So we learn the function

(*) But if we change $u_0(x)$?

(#) The solution map is way more useful!

ex: $u_0(x) \in L^1(\mathbb{R})$

Then

$$u(x,t) = (G_t * u_0)(x)$$

$$= \int_{\mathbb{R}} G(x-y, t) \, u_0(y) \, dy$$

$$G_t = \frac{1}{\sqrt{4\pi t}} \, e^{-x^2/4t}$$

What about forcing?

$$u_t = u_{xx} + f$$

$$u = G_t * u_0$$
$$+ \int_0^t \int_{\mathbb{R}} G(x-y, t-s) \, f(y, s) \, dy \, ds$$

So much more useful
than just one function!

(*) This shift from functions to operators is the key idea here

(*) So how do we approach learning an operator?

- Take inspiration from traditional deep learning

- Consider an NN layer

$$v_i^{l+1} = \sigma\left(\sum_{j=1}^{n} W_{ij} v_j^{l} + b_i\right)$$

activation
$\neq$
ReLu

weights

bias

what if we take $n \to \infty$
for $v^l = \langle v_1^l, v_2^l, \ldots, v_n^l \rangle$?

$$\sum_{j=1}^{n} W_{ij} v_j^l \longrightarrow \int k(x_i, y) v^l(y) \, dy$$

(*) Dense layers are discretized integral operators

(*) A neural operator in this setting is just a NN with layer

$$v^{l+1}(x) = \sigma \left( A_\theta v^l(x) + \int_\Omega \kappa_\theta(x,y) v^l(y) \, dy \right)$$

(by analogy)

(*) But now that we are on function spaces.... we can discretize back!

$$v^{l+1}(x_i) = \sigma \left( A_\theta v^l(x_i) + \sum_{j=1}^{n} \kappa_\theta(x_i, x_j) \, v^l(x_j) \, w_j \right)$$

↑ quadrature weights

(*) If the PDE grid is irregular use a graph!

- Nodes = points $x_i$
- Edges connect $x_i$ to nearby $x_j$

(∗) $n_\theta(x_i, x_j)$ plays the role of an edge weight in this local "message-passing" scheme

(∗) If we have a regular grid, then we can do better.

Especially with translation invariance

$$V^{l+1} = \sigma \left( A_\theta z^l + \eta_\theta * v^l \right)$$

you know what's coming

$$\overbrace{\eta_\theta * v^l} = \hat{\eta}_\theta \hat{v}^l$$

Building the Fourier piece has 3 steps.

(i) $\hat{v}^l(k) = \mathcal{F}\{v^l\}(k)$

(ii) $\hat{v}^{l+1}(k) = R_\theta(k) \hat{v}^l(k)$

⊕ And this is what you learn

⊛ Sometimes diagonal, Sometimes
a low pass filter...
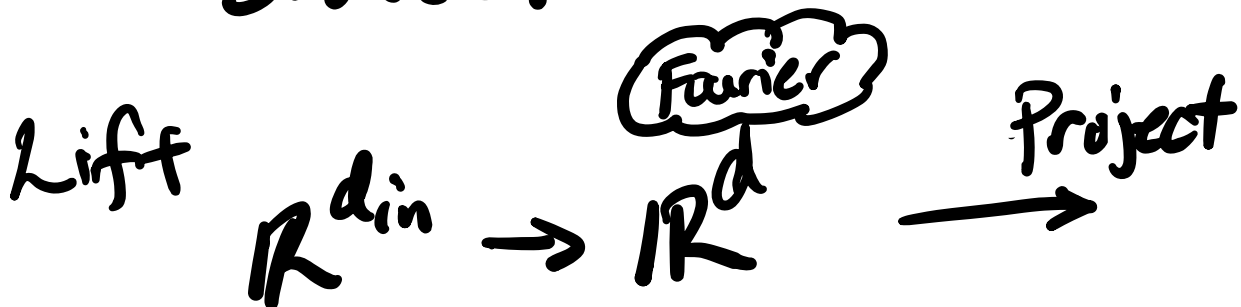(#) Signal processing people
valuable here

$$(l^{ric}) \quad \eta_0 * \nu^l$$

$$= \mathcal{F}^{-1} \{ \hat{\eta_0} \hat{\nu}^l \} (x)$$

$$= \mathcal{F}^{-1} \{ R_0 \hat{\nu}^l \} (x)$$

$$\nu^{l+1} = \sigma \left( \Lambda_0 \nu^l + \mathcal{F}^{-1} \{ R_0 \hat{\nu}^l \} \right)$$

(*) This is just one
layer

(**) What's the whole
structure?

Lift
$$\mathbb{R}^{d_{in}} \xrightarrow{\text{Fourier}} \mathbb{R}^d \xrightarrow{\text{Project}}$$

$$v^0(x) = P_{lift}\left(\underset{\underset{\text{input function}}{\uparrow}}{a(x)}\right)$$

$$v^{l+1} = \ldots \quad (\text{insert here})$$

$$l = 0, 1, \ldots, L-1$$

$$u(x) = P_{proj}\left(v^L(x)\right)$$
$$(\text{output decoding})$$

$$G_\theta(a) = P_{proj} \circ$$

$$\left(\lambda_\theta^{(L)} \circ \ldots \circ \lambda_\theta^{(1)}\right) \circ P_{lift}(a)$$

(!) This, the paper

claims, is the **only** resolution invariant neural operator that universally approximates

(✳) The paper trashes Deep ONets (branch/trunk)

Let's pose the learning framework....

- Let $A$ and $U$ be Banach spaces

of functions on

$$D \subset \mathbb{R}^d, \quad D' \subset \mathbb{R}^{d'}$$

- Denote the true map as

$$G^\dagger : A \to \mathcal{U}$$

- Define the Bochner norm

$$\|f\|^2_{L^2_\mu(A;\mathcal{U})} = \int_A \|f(a)\|^2_{\mathcal{U}} \, d\mu(a)$$

$$= \mathop{\mathbb{E}}_{a \sim \mu} \left\{ \|f(a)\|^2_{\mathcal{U}} \right\}$$

Theoretical: $\min\limits_{\theta \in \mathbb{R}^P} \|G^\dagger - G_\theta\|^2_{L^2_\mu(A;\mathcal{U})}$

Practical: $\min\limits_{\theta \in \mathbb{R}^P} \frac{1}{N} \sum\limits_{i=1}^{N} \|u^{(i)} - G_\theta(a^{(i)})\|^2_{\mathcal{U}}$

(a) Define the discretization error as

$$|J_{theoretical} - J_{practical}|$$

Assume

$*$ $\sup_x |\int g - \Sigma g| \leq Ch^q$

$*$ $\eta_{\theta,h}$ lipschitz in $h$

$*$ $||S_h G^+(a) - G_\theta(S_h a)||^2_{h,\mu} \leq \varepsilon^2$

then

$$||S_{\tilde{h}} G^+(a) - G_\theta(S_{\tilde{h}} a)||_{h,\mu}$$

$$\leq \varepsilon + \hat{C}(h^q + \tilde{h}^q)$$

approximation error       training & test
                         quadrature error

# Universal Approximation

$$\exists \ L \ , \ d \ , \ \theta$$

depth ; width , parameters

s.t.

$$\sup_{a \in K} \| G^T(a) - G_\theta(a) \|_{u} < \varepsilon$$

where $K \subset A$

compact

$$G^T : A \to u$$

Continuous (with respect to the norm topology on $A$ ; $u$ )