

Generalization Bounds for Solving Navier-Stokes PDE by Neural Nets

Dibyakanti Kumar, Department of Computer Science, The University of Manchester

Collaborators

This work is with,

Sébastien André-Sloan @ UManchester, (PhD Candidate, CSE),

Alejandro Frangi @ UManchester (Computer Science)

&

Anirbit Mukherjee @ UManchester (Computer Science)

Physics Informed Neural Networks

Physics Informed Neural Nets (PINNs)

Let $u(x, t)$ be the actual solution

PDE

$$u_t + \mathcal{N}_x[u] = 0, \quad x \in D, t \in [0, T]$$

Initial Condition

$$u(x, 0) = h(x), \quad x \in D$$

Boundary Conditions

$$u(x, t) = g(x, t), \quad t \in [0, T], x \in \partial D$$

Physics Informed Neural Nets (PINNs)

Let the neural solution be $u_\theta(x, t)$

Functional Residual

$$\mathcal{R}_{pde,\theta}(x, t) := \partial_t u_\theta + \mathcal{N}_x[u_\theta(x, t)]$$

Conditional Residual (I.C.)

$$\mathcal{R}_{t,\theta}(x) := u_\theta(x, 0) - h(x)$$

Conditional Residual (B.C.)

$$\mathcal{R}_{b,\theta}(x, t) := u_\theta(x, t) - g(x, t)$$

Physics Informed Neural Nets (PINNs)

Let the neural solution be $u_\theta(x, t)$

Functional Residual

$$\mathcal{R}_{pde,\theta}(x, t) := \partial_t u_\theta + \mathcal{N}_x[u_\theta(x, t)]$$

Conditional Residual (I.C.)

$$\mathcal{R}_{t,\theta}(x) := u_\theta(x, 0) - h(x)$$

Conditional Residual (B.C.)

$$\mathcal{R}_{b,\theta}(x, t) := u_\theta(x, t) - g(x, t)$$

PINN Loss Function

$$\mathcal{L}_S(\theta) := \frac{1}{N_{pde}} \sum_{i=1}^{N_{pde}} \mathcal{R}_{pde,\theta}(x_r^i, t_r^i)^2 + \lambda_1 \frac{1}{N_t} \sum_{i=1}^{N_t} \mathcal{R}_{t,\theta}(x_t^i)^2 + \lambda_2 \frac{1}{N_b} \sum_{i=1}^{N_b} \mathcal{R}_{b,\theta}(x_b^i, t_b^i)^2$$

Why Physics Informed Neural Nets (PINNs) ?

$$\mathcal{E}_G := \|\mathbf{u} - \mathbf{u}_\theta\|$$

$$\tilde{\mathcal{E}}_T := \left(\sum_{n=1}^N w_n |\mathcal{R}_\theta|^p \right)^{1/p}$$

Risk Upperbound (Theorem 2.6)¹

$$\mathcal{E}_G \leq C_{pde} \tilde{\mathcal{E}}_T + C_{pde} C_{quad}^{\frac{1}{p}} N^{-\frac{\alpha}{p}}$$

It was proven in the above mentioned theorem – for the first time – that one can minimize certain empirical errors to find provably good approximations to any PDE.

¹Mishra and Molinaro, “Estimates on the generalization error of physics-informed neural networks for approximating PDEs”.

Why Physics Informed Neural Nets (PINNs) ?

Upperbounds of L2-risk and Residual for any PDE (Theorem 3.5)²

$$\|u - u_\theta\|_{L^2} + \|R_{pde,\theta}\|_{L^2} \leq C_1 \cdot \text{poly}(d) \cdot C_u$$

- This showed that there **exists** neural networks that can simultaneously reduce the PDE loss and the L2-risk.
- **Does not provide a finite-sample bound.**

²De Ryck and Mishra, “Generic bounds on the approximation error for physics-informed (and) operator learning”.

Why Physics Informed Neural Nets (PINNs) ?

Upperbounds of Risk for Linear PDE (Theorem 2)³

$$\|u - u_\theta\|_{H^s} \leq \sqrt{C_u(\mathbb{E}[u_\theta] - \inf_\psi \mathbb{E}[u_\psi]) + C_{pde} \inf_\psi \|u_\psi - u\|}$$

- This showed that there **exists** neural networks that can reduce the risk in Sobolev norm for PINNs.
- **Does not provide a finite-sample bound.**

³Zeinhofer, Masri, and Mardal, “A unified framework for the error analysis of physics-informed neural networks”.

Why Physics Informed Neural Nets (PINNs) ?

Upperbounds of L2-risk for NS (Theorem 3.4)⁴

$$\|\mathbf{u} - \mathbf{u}_\theta\|_{L^2} \leq C_2 \cdot C_{\nabla u}$$

- This showed that there **exists** neural networks that can reduce the L2-risk for PINNs when solved for Navier-Stokes.
- Does not provide a finite-sample bound.

⁴De Ryck, Jagtap, and Mishra, “Error estimates for physics-informed neural networks approximating the Navier–Stokes equations”.

Why PINNs for Navier–Stokes?

Recent studies have demonstrated the effectiveness of solving fluid flow problems with a physics-informed loss.

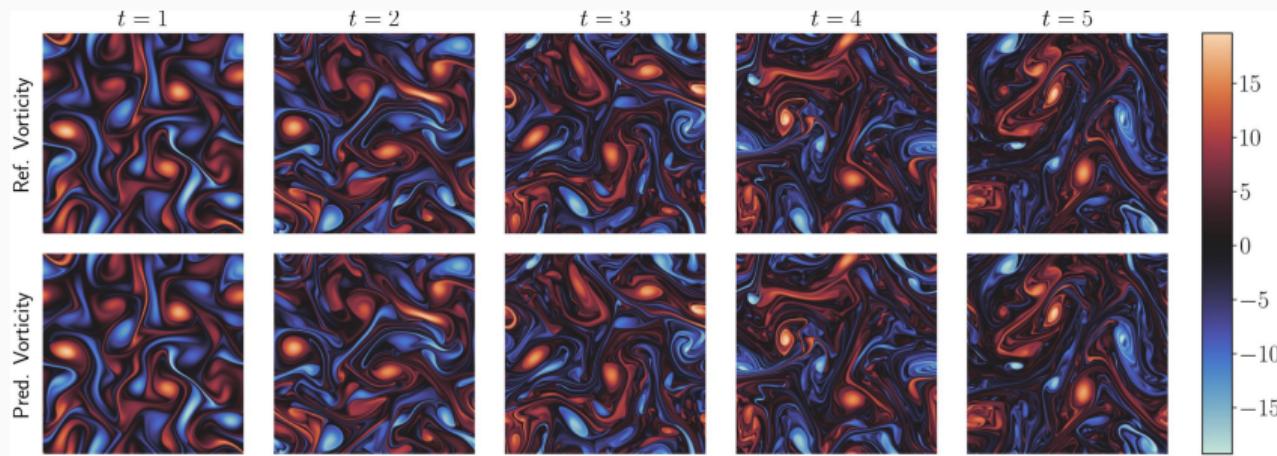


Figure 1: Time evolution of the 2D Kolmogorov flow solved using PirateNet, as reported in a recent study⁵.

⁵Wang et al., “Simulating three-dimensional turbulence with physics-informed neural networks”.

What are We Missing?

But “learning” is far more than just an approximation guarantee!

We need to explain why the PINNs gives good approximations at
points that never occurred in the training data.

What are We Missing?

But “learning” is far more than just an approximation guarantee!

We need to explain why the PINNs gives good approximations at
points that never occurred in the training data.

This brings us to the core question in learning theory
- that of quantifying “**generalization error**”.

Rademacher Complexity

What Makes Up a Deep Learning Model?

The following choices must be made with care:

1. Network architecture (e.g., size, activation functions, ...)
2. Training algorithm
3. Data
4. Loss Function

What Makes Up a Deep Learning Model?

The following choices must be made with care:

1. Network architecture (e.g., size, activation functions, ...)
2. Training algorithm
3. Data
4. Loss Function

What can we expect to achieve in theory?

- The **optimal combination** of these four components
- A **provably good** choice of these components

What Makes Up a Deep Learning Model?

The following choices must be made with care:

1. Network architecture (e.g., size, activation functions, ...)
2. Training algorithm
3. Data
4. Loss Function

What can we expect to achieve in theory?

- The **optimal combination** of these four components
- A **provably good** choice of these components ← This is where **Rademacher complexity** can guide us.

What is Rademacher Complexity?

Suppose being given a function class \mathcal{H} mapping into \mathbb{R} , $\ell(\cdot)$ be a L -Lipschitz loss function and a m -sized data-set \mathcal{D} of points $\{x_i \mid i = 1, \dots, m\}$ in domain (\mathcal{H}) .

Then, for $\epsilon_i \sim \text{Uniform}(\pm 1)$, the **empirical Rademacher complexity** is,

$$\hat{\mathcal{R}}_m(\ell(\mathcal{H})) = \mathbb{E}_\epsilon \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \epsilon_i \ell(h(x_i)) \right].$$

If the data \mathcal{D} above are sampled from a distribution P , then the **average Rademacher complexity** is,

$$\mathcal{R}_m(\ell(\mathcal{H})) = \mathbb{E}_{x_1, \dots, x_m \sim P} \left(\mathbb{E}_\epsilon \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \epsilon_i \ell(h(x_i)) \right] \right).$$

Is Rademacher Inevitable?



I am inevitable.

Don't worry — the irony isn't lost on me.

Is Rademacher Inevitable?

Theorem

$$\mathbb{E}_{x_1, \dots, x_m \sim P} \left[\sup_{h \in \mathcal{H}} \left(\frac{1}{m} \cdot \sum_{i=1}^m \ell(h(x_i)) - \mathbb{E}_{x \sim P} [\ell(h(x))] \right) \right] \leq 2 \cdot \text{Rad}_m(\ell(\mathcal{H}))$$

- In real-world uses ℓ is the loss function of the predictor h – and hence the LHS measures the “data averaged worst gap possible between the empirical loss and the true risk of the predictor”.

Is Rademacher Inevitable?

Theorem

$$\mathbb{E}_{x_1, \dots, x_m \sim P} \left[\sup_{h \in \mathcal{H}} \left(\frac{1}{m} \cdot \sum_{i=1}^m \ell(h(x_i)) - \mathbb{E}_{x \sim P} [\ell(h(x))] \right) \right] \leq 2 \cdot \text{Rad}_m(\ell(\mathcal{H}))$$

- In real-world uses ℓ is the loss function of the predictor h – and hence the LHS measures the “data averaged worst gap possible between the empirical loss and the true risk of the predictor”.
- This philosophically says that the ability of the predictor class to replicate uniform binary noise on a discrete set (its “capacity”) has to be small for it to have the ability to “learn”.

Gen. Error Bound for Linear 2nd-Order PDEs (Theorem 3.2)⁶

$$\mathbb{E} \left[\sup_{\theta} (\mathcal{L}_{\mathcal{S}}(\theta) - \mathbb{E}_{\mathcal{S}} \mathcal{L}_{\mathcal{S}}(\theta)) \right] \leq \frac{C_1}{N_r \sqrt{N_r}} + \frac{C_2 \log N_r}{\sqrt{N_r}} + \frac{C_3}{N_b \sqrt{N_b}} + \frac{C_4 \log N_b}{\sqrt{N_b}}$$

- This is a finite-sample complexity bound to the generalization error.

⁶Hu et al., "When Do Extended Physics-Informed Neural Networks (XPINNs) Improve Generalization?"

Rademacher Complexity for Navier-Stokes

Rademacher Complexity for (d+1)-Navier-Stokes

Let the output of the neural net we aim to train be $\mathbf{u}_\theta(\mathbf{x}, t) \in \mathbb{R}^d$ and $p_\theta(\mathbf{x}, t) \in \mathbb{R}$.

Then the residual terms for this predictor are,

Conditional Residual

$$\mathcal{R}_{t,\theta} := \mathbf{u}_\theta(t = t_0) - \mathbf{u}(t = t_0)$$

The residual term for the functional loss is

Functional Residual

$$\mathcal{R}_{\text{pde},\theta} := \partial_t \mathbf{u}_\theta + (\mathbf{u}_\theta \cdot \nabla) \mathbf{u}_\theta + \nabla p - \nu \nabla^2 \mathbf{u}_\theta; \quad \mathcal{R}_{\text{div},\theta} := \nabla \cdot \mathbf{u}_\theta$$

Generalization Error

$$\begin{aligned}\mathcal{E}_{gen,\theta} := & \frac{1}{N_r} \sum_{i=1}^{N_r} \mathcal{R}_{pde,\theta}(\mathbf{x}_{ri}, t_{ri}) + \frac{\lambda_0}{N_r} \sum_{i=1}^{N_r} \mathcal{R}_{div,\theta}(\mathbf{x}_{ri}, t_{ri}) + \frac{\lambda_1}{N_0} \sum_{j=1}^{N_0} \mathcal{R}_{t,\theta}(\mathbf{x}_{0j}) \\ & - (\mathbb{E}_{(\mathbf{x}_r, t_r) \sim \mathcal{D}_1} [\mathcal{R}_{pde,\theta}(\mathbf{x}_r, t_r) + \lambda_0 \mathcal{R}_{div,\theta}(\mathbf{x}_r, t_r)] + \lambda_1 \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{D}_2} [\mathcal{R}_{t,\theta}(\mathbf{x}_0)])\end{aligned}$$

Generalization Error

$$\begin{aligned}\mathcal{E}_{gen,\theta} := & \frac{1}{N_r} \sum_{i=1}^{N_r} \mathcal{R}_{pde,\theta}(\mathbf{x}_{ri}, t_{ri}) + \frac{\lambda_0}{N_r} \sum_{i=1}^{N_r} \mathcal{R}_{div,\theta}(\mathbf{x}_{ri}, t_{ri}) + \frac{\lambda_1}{N_0} \sum_{j=1}^{N_0} \mathcal{R}_{t,\theta}(\mathbf{x}_{0j}) \\ & - (\mathbb{E}_{(\mathbf{x}_r, t_r) \sim \mathcal{D}_1} [\mathcal{R}_{pde,\theta}(\mathbf{x}_r, t_r) + \lambda_0 \mathcal{R}_{div,\theta}(\mathbf{x}_r, t_r)] + \lambda_1 \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{D}_2} [\mathcal{R}_{t,\theta}(\mathbf{x}_0)])\end{aligned}$$

Theorem 1

The worst case data-averaged generalization error can be bounded as,

$$\mathbb{E} \left[\sup_{\theta} (\mathcal{E}_{gen,\theta}) \right] \leq \frac{C_r}{\sqrt{N_r}} + \frac{C_0}{\sqrt{N_0}}.$$

Theorem 1

The worst case data-averaged generalization error can be bounded as,

$$\mathbb{E} \left[\sup_{\theta} (\mathcal{E}_{gen, \theta}) \right] \leq \frac{C_r}{\sqrt{N_r}} + \frac{C_0}{\sqrt{N_0}}.$$

- The loss is Lipschitz (i.e. Huber loss)
- Depth-2 neural nets
- tanh activation

Rademacher Complexity for (d+1)-Navier-Stokes

Theorem 1

The worst case data-averaged generalization error can be bounded as,

$$\mathbb{E} \left[\sup_{\theta} (\mathcal{E}_{gen, \theta}) \right] \leq \frac{C_r}{\sqrt{N_r}} + \frac{C_0}{\sqrt{N_0}}.$$

- The loss is Lipschitz (i.e. Huber loss)
- Depth-2 neural nets
- tanh activation

This is possibly the first generalization error bound for solving a non-linear PDE by neural nets.

An Overview of the Main Lemmas

Lipschitz Composition in Rademacher Complexity

Talagrand's Contraction Lemma

If ℓ is an L_ℓ -Lipschitz loss function, then

$$\mathcal{R}_m(\ell(\mathcal{H})) \leq L_\ell \cdot \mathcal{R}_m(\mathcal{H})$$

Why Do We Need a Contraction Lemma?

The key idea behind the proof is to establish a suitable **contraction lemma** for each component of the loss function.

After applying these contraction lemmas, we aim to reduce the problem to analyzing the **Rademacher complexity of a linear model**.

Why Do We Need a Contraction Lemma?

The key idea behind the proof is to establish a suitable **contraction lemma** for each component of the loss function.

After applying these contraction lemmas, we aim to reduce the problem to analyzing the **Rademacher complexity of a linear model**.

In our case, we developed **two contraction lemmas**: one for the non-linear term, and one for the remaining terms.

Contraction Lemma I

- $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is L_ϕ -Lipschitz
- For any $\epsilon \in \{\pm 1\}^N$, $\sup_{\mathbf{w}} \langle \epsilon, \phi(\mathbf{Z}^\top \mathbf{w}) \rangle \geq 0$, where $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$
- $\sum_{m=1}^p |f(\mathbf{W})_m| \leq B$, where $f(\mathbf{W}) : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^p$

Contraction Lemma I

$$\mathbb{E}_{\epsilon \sim \{\pm 1\}^N} \left[\sup_{\mathbf{W}} \left(\frac{1}{N} \sum_{i=1}^N \epsilon_i \langle f(\mathbf{W}), \phi(\mathbf{W} \mathbf{z}_i) \rangle \right) \right] \leq \frac{2BL_\phi}{N} \mathbb{E}_{\epsilon \sim \{\pm 1\}^N} \left[\sup_{\bar{\mathbf{w}}} \sum_{i=1}^N \epsilon_i \langle \bar{\mathbf{w}}, \mathbf{z}_i \rangle \right]$$

Contraction Lemma II : Non-Linear Term

- $\phi_1, \phi_2 : \mathbb{R} \rightarrow \mathbb{R}$ are L_{ϕ_1}, L_{ϕ_2} -Lipschitz, and $\phi_1 \leq B_{\phi_1}$ and $\phi_2 \leq B_{\phi_2}$
- $\phi_1(\cdot)\phi_2(\cdot)$ is an odd function
- $f_m(\mathbf{W}) : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^p$ such that $\forall m \in [d], \sum_{q_1, q_2=1}^p |f_m(\mathbf{W})_{q_1}| |\mathbf{a}_{q_2}| \leq B$

Contraction Lemma II

$$\begin{aligned} & \mathbb{E}_{\epsilon \sim \{\pm 1\}^N} \left[\sup_{\mathbf{W}} \left(\frac{1}{N} \sum_{i=1}^N \sum_{m=1}^d \epsilon_i \langle f_m(\mathbf{W}), \phi_1(\mathbf{W}\mathbf{z}_i) \rangle \langle \mathbf{a}_m, \phi_2(\mathbf{W}\mathbf{z}_i) \rangle \right) \right] \\ & \leq \frac{2B(B_{\phi_1}L_{\phi_2} + B_{\phi_2}L_{\phi_1})}{N} \mathbb{E}_{\epsilon \sim \{\pm 1\}^N} \left[\sup_{\bar{\mathbf{w}}} \left(\sum_{i=1}^N \epsilon_i \langle \bar{\mathbf{w}}, \mathbf{z}_i \rangle \right) \right] \end{aligned}$$

Putting It All Together

$$\mathbb{E} \left[\sup_{\theta} (\mathcal{E}_{gen, \theta}) \right] \leq \frac{C_r}{\sqrt{N_r}} + \frac{C_0}{\sqrt{N_0}}.$$

$$C_r := L_\ell B_w C_z (2B_{f_1} L_{\sigma'} + 2B_{f_2} (B_\sigma L_{\sigma'} + B_{\sigma'} L_\sigma) + 2B_{f_3} L_{\sigma'} + 4\nu B_{f_4} (L_{\sigma'} + L_\sigma) + 2\lambda_0 B_{f_5} L_{\sigma'})$$

$$C_0 := 2\lambda_1 L_\ell B_w C_{x_0} L_\sigma B_a$$

Looking Ahead

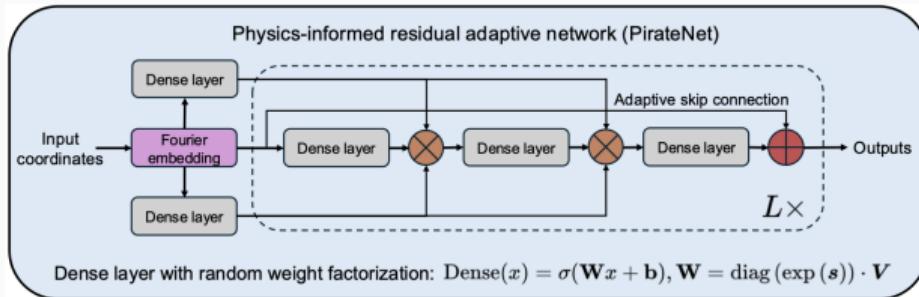


Figure 2: PirateNet architecture, as used in a recent study⁷.

SOTA architectures like PirateNet are complex, and more theory is needed to explain their success on Navier-Stokes.

Much work remains to move from provably good settings to truly optimal ones that also consider loss optimizability.

⁷Wang et al., "Simulating three-dimensional turbulence with physics-informed neural networks".

Thank You!

To learn more about our work, scan the QR code or visit:
<https://dibyakanti.github.io/publication/>

