

Mean field limit in inverse problem

Yutian He, Joseph Hunter, Vasily Ilin, Ian McPherson, Kouakou
Innocent NDRI, Yantao Wu, Jaeyoung Yoon

SLMath: Particle interactive systems

28th June 2024

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 2/28

Problem Setup: Mean-field version

Function $u_\rho(t, x) : [0, 1] \times \mathbb{R} \rightarrow \mathbb{R}$ given by the dynamics

$$\begin{cases} \frac{d}{dt} u_\rho(t, x) = \int_{\mathbb{R}} f(u_\rho(t, x), \theta) \rho(\theta) d\theta, & t \in [0, 1] \\ u_\rho(t = 0, x) = x \end{cases} \quad (1)$$

1. We are given the “activation function” f and data $D(x) = u_{\rho_*}(t = 1, x)$.
2. Want to approximate the distribution of “weights” $\rho_*(\theta)$.

Solution: Minimizing loss functional

$$E[\rho] = \frac{1}{2} \int_{\Omega_x} |u_\rho(1, x) - D(x)|^2 d\mu(x)$$

Problem Setup: Particle version

Function $u_{\rho_N} : [0, 1] \times \mathbb{R} \rightarrow \mathbb{R}$ has the following dynamics

$$\begin{cases} \frac{d}{dt} u_{\rho_N} = \frac{1}{N} \sum_{i=1}^N f(u_{\rho_N}, \theta_i), & t \in [0, 1] \\ u_{\rho_N}(t = 0, x) = x \end{cases} \quad (2)$$

Use the particle approximation

$$\rho_{N,s} = \frac{1}{N} \sum_i \delta_{\theta_i(s)}$$

$$\frac{d}{ds} \theta_i(s) = -\partial_{\theta_i} E[\rho_{N,s}] = W_{\rho_{N,s}}(\theta_i)$$

$$= - \int_0^1 \int v_{\rho_{N,s}}(t, x) \partial_{\theta_i} f(u_{\rho_{N,s}}, \theta) dx dt, \quad v = \text{adjoint term}$$

$$\frac{d}{dt} v_{\rho_{N,s}} = -v_{\rho_{N,s}} \int_{\mathbb{R}} \frac{\partial f}{\partial u}(u_{\rho_{N,s}}, \theta) \rho_{N,s}(\theta) d\theta, \quad v_{\rho_{N,s}}(t = 1) = D - u_{\rho_{N,s}}(t = 1).$$

Interpretation via Neural ODE

Dynamics (2) describes a Neural Network:

1. Layer $t \in [0, 1]$ and input x .
2. $u_{\rho_N}(t, x) =$ output from t -th layer. First layer is identity $u(t = 0, x) = x$.
3. Minimize sum of squares $E[\rho_N] = \int |D(x) - u_{\rho_N}(t = 1, x)|^2 d\mu(x)$.
4. Residual network architecture, finite width N , infinitely deep, layer-wise homogeneous neural network.

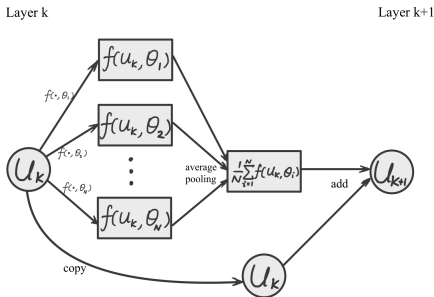


Figure: One layer of NN

Gradient Flows

Wasserstein Gradient Flow

$$\begin{cases} \partial_s \rho_s = -\nabla_{\mathcal{W}_2} E[\rho_s] = \partial_\theta \left(\rho_s \partial_\theta \frac{\delta E}{\delta \rho} \Big|_{\rho_s} \right), & s \in [0, T] \\ \rho_{s=0} = \rho_0 \end{cases} \quad (3)$$

Euclidean Gradient Flow

$$\begin{cases} \partial_s \Theta_s = -\nabla_\Theta E[\rho_{N,s}] & s \in [0, T] \\ \Theta_{s=0} = \Theta_0 \end{cases} \quad (4)$$

Main Goal: Compare how empirical measure

$$\rho_{N,s} = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_{i,s}}$$

approximates ρ_s as $N \rightarrow \infty$.

Remark: convergence as $s \rightarrow \infty$ is a different story.

Main Theorem

Theorem (Mean-field limit)

Assume the following finiteness conditions on f and ρ_s :

$$C_f^0 = \|\partial_\theta f(y, \theta)\|_{L^\infty(y, \theta)} + \|\partial_y f(y, \theta)\|_{L^\infty(y, \theta)} < \infty$$

$$C_f^1 = 1 + \|\partial_\theta^2 f(y, \theta)\|_{L^\infty(y, \theta)} + \|\partial_{\theta, y}^2 f(y, \theta)\|_{L^\infty(y, \theta)} + \|\partial_y^2 f(y, \theta)\|_{L^\infty(y, \theta)} < \infty$$

$$\mathbb{E}(E[\rho_s] + E[\rho_{N,s}] + E[\bar{\rho}_{N,s}]) < \infty, \quad E[\rho] = \int |D(x) - u_{\rho_s}(t=1, x)|^2 d\mu(x)$$

$$\int |\theta|^5 d\rho_s(\theta) < \infty. \quad \text{Then, for all } s \in [0, T]$$

$$\mathbb{E}\mathcal{W}_2(\rho_s, \rho_{N,s}) \leq \frac{C}{N^{1/4}} + \frac{C''}{N^{1/2}} \exp(C's) = \mathcal{O}(N^{-1/4}),$$

$$\mathbb{E}\mathcal{W}_1(\rho_s, \rho_{N,s}) = \mathcal{O}(N^{-1/2})$$

$$C = \text{universal}, \quad C' = 16 (C_f^1)^2 e^{5C_f^0} \mathbb{E} \left(E[\rho_s]^{\frac{1}{2}} + E[\rho_{N,s}]^{\frac{1}{2}} + E[\bar{\rho}_{N,s}]^{\frac{1}{2}} \right) + 1$$

$$C'' = C' \|f(u_\rho, \theta)\|_{L^4(\mu(x) \otimes \rho(x))}$$

Proof sketch

$$\begin{aligned}
 \frac{d}{ds} |\theta_i - \bar{\theta}_i| &= |W_{\hat{\rho}_N}(\theta_i) - W_{\rho_s}(\bar{\theta}_i)| \\
 &\leq |W_{\hat{\rho}_N}(\theta_i) - W_{\hat{\rho}_N}(\bar{\theta}_i)| + |W_{\hat{\rho}_N}(\bar{\theta}_i) - W_{\bar{\rho}_N}(\bar{\theta}_i)| + |W_{\bar{\rho}_N}(\bar{\theta}_i) - W_{\rho_s}(\bar{\theta}_i)| \\
 &\leq L_A |\theta_i - \bar{\theta}_i| + L_B \left(\frac{1}{N} \sum_{j=1}^N |\theta_j - \bar{\theta}_j| \right) + L_C \frac{1}{\sqrt{N}}
 \end{aligned}$$

Sum over i and combine Lipschitz constants L_A, L_B, L_C :

$$\frac{d}{ds} \left(\frac{1}{N} \sum_{i=1}^N |\theta_{i,s} - \bar{\theta}_{i,s}|^2 \right) \leq \frac{(C'')^2}{N} + C' \left(\frac{1}{N} \sum_{i=1}^N |\theta_{i,s} - \bar{\theta}_{i,s}|^2 \right)$$

By Grönwall's inequality,

$$\mathcal{W}_2^2(\bar{\rho}_{N,s}, \rho_{N,s}) \leq \frac{1}{N} \sum_{i=1}^N |\theta_{i,s} - \bar{\theta}_{i,s}|^2 \leq \frac{(C'')^2}{N} \exp(C's)$$

$$\mathcal{W}_2(\rho_s, \rho_{N,s}) \leq \mathcal{W}_2(\rho_s, \bar{\rho}_{N,s}) + \mathcal{W}_2(\bar{\rho}_{N,s}, \rho_{N,s}) \leq \frac{C}{N^{1/4}} + \frac{C''}{N^{1/2}} \exp(C's)$$

Lipschitz lemmas

Lemma (Flow Velocity $W_\rho(\theta)$ is Lipschitz in θ)

$$|W_{\rho_{N,s}}(\theta_1) - W_{\rho_{N,s}}(\theta_2)| \leq 2C_f^1 e^{C_f^0} E[\rho_{N,s}]^{\frac{1}{2}} |\theta_1 - \theta_2|.$$

Lemma (Flow Velocity $W_\rho(\theta)$ is Lipschitz-Like in $\hat{\rho}_N, \bar{\rho}_N$)

$$|W_{\rho_{N,s}}(\theta) - W_{\bar{\rho}_{N,s}}(\theta)| \leq 8 (C_f^1)^2 e^{5C_f^0} \left(E[\rho_{N,s}]^{\frac{1}{2}} + E[\bar{\rho}_{N,s}]^{\frac{1}{2}} \right) \left(\frac{1}{N} \sum_{i=1}^N |\theta_i - \bar{\theta}_i| \right)$$

Lemma (Flow Velocity $W_\rho(\theta)$ is Lipschitz-Like in $\rho, \bar{\rho}_N$)

If $\|f(u_\rho, \theta)\|_{L^4(\mu(x) \otimes \rho(x))} < \infty$, then

$$\mathbb{E}_\Omega |W_\rho(\theta) - W_{\bar{\rho}_N}(\theta)| \leq \frac{8}{\sqrt{N}} C_f^1 e^{5C_f^0} \left(E[\rho]^{\frac{1}{2}} + E[\bar{\rho}_N]^{\frac{1}{2}} \right) \|f(u_\rho, \theta)\|_{L^4(\mu(x) \otimes \rho(x))}$$

Numeric analysis

Minimization problem

In numeric, we can not have a full information for a distribution ρ . Instead, we approximate $\rho(\theta)$ with $\rho_N = \frac{1}{N} \sum_i \delta_{\theta_i}$ and solve the following minimization problem:

$$\min_{\Theta_N \in \mathbb{R}^N} E(\Theta_N) = \min_{\Theta_N \in \mathbb{R}^N} \frac{1}{2} \int_{\Omega_x} |u_{\rho_N}(1, x) - D(x)|^2 dx,$$

where the dynamic $u_{\rho_N}(t, x)$ follows

$$\begin{cases} \frac{d}{dt} u_{\rho_N}(t, x) = \frac{1}{N} \sum_{i=1}^N f(u_{\rho_N}(t, x), \theta_i), & \forall t \in (0, 1), \\ u_{\rho_N}(0, x) = x, & \forall x \in \Omega_x. \end{cases}$$

Lagrange Multiplier

Since it is hard to compute $\nabla_{\Theta_N} E(\Theta_N)$ directly, we define a Lagrangian functional \mathcal{L} by

$$\begin{aligned} \mathcal{L}(u, \Theta_N, \eta, \tilde{\eta}) := & \frac{1}{2} \int_{\Omega_x} |u(1, x) - D(x)|^2 dx - \int_{\Omega_x} (u(0, x) - x) \tilde{\eta}(x) dx \\ & - \int_0^1 \int_{\Omega_x} \left(\frac{d}{dt} u(t, x) - \frac{1}{N} \sum_{i=1}^N f(u(t, x), \theta_i) \right) \eta(t, x) dx dt, \end{aligned}$$

where $\eta(t, x)$ and $\tilde{\eta}(x)$ are Lagrange multipliers. Note that the first-order optimality condition is

$$\frac{\delta \mathcal{L}}{\delta \eta} = \frac{\delta \mathcal{L}}{\delta \tilde{\eta}} = \frac{\delta \mathcal{L}}{\delta u(t, x)} = \frac{\delta \mathcal{L}}{\delta u(0, x)} = \frac{\delta \mathcal{L}}{\delta u(1, x)} = \frac{\partial \mathcal{L}}{\partial \Theta_N} = 0.$$

Strategy

While keeping

$$\frac{\delta \mathcal{L}}{\delta \eta} = \frac{\delta \mathcal{L}}{\delta \tilde{\eta}} = \frac{\delta \mathcal{L}}{\delta u(t, x)} = \frac{\delta \mathcal{L}}{\delta u(0, x)} = \frac{\delta \mathcal{L}}{\delta u(1, x)} = 0, \quad (5)$$

make a gradient flow for $\Theta_N(s)$ as

$$\frac{d}{ds} \Theta_N(s) = - \frac{\partial \mathcal{L}}{\partial \Theta_N} = - \frac{1}{N} \sum_{i=1}^N \int_0^1 \int_{\Omega_x} \nabla_{\Theta} f(u(t, x), \theta_i(s)) \eta(t, x) dx dt.$$

The conditions (5) gives the **forward system** (original constraint) and **adjoint system**:

$$(F) \quad \begin{cases} \partial_t u(t, x) = \frac{1}{N} \sum_{i=1}^N f(u(t, x), \theta_i), \\ u(0, x) = x, \quad \forall (t, x) \in (0, 1) \times \Omega_x. \end{cases}$$

Algorithm at a glance

1. Flow for Θ in pseudo time s :

$$\Theta_N(s_{k+1}) = \Theta_N(s_k) - \frac{h_s h_t}{NM} \sum_{i=1}^N \sum_{l=1}^{N_t} \sum_{j=1}^M \nabla_{\Theta} f(u(t_l, x_j), \theta_i) \eta(t_l, x_j).$$

2. Forward system in physical time t :

$$u(t_{l+1}, x) = u(t_l, x) + \frac{h_t}{N} \sum_{i=1}^N f(u(t_l, x), \theta_i), \quad u(0, x) = x,$$

3. Adjoint system in physical time t :

$$\begin{aligned} \tilde{\eta}(t_{l+1}, x) &= \tilde{\eta}(t_l, x) + \frac{h_t}{N} \sum_{i=1}^N \partial_u f(u((N_t - l)h_t, x), \theta_i), \\ \tilde{\eta}(0, x) &= u(1, x) - D(x). \end{aligned}$$

- Uniform time sequences $\{t_l\}_{l=0}^{N_t}$ and $\{s_k\}_{k=0}^{N_s}$ satisfy

$$0 = t_0 < t_1 < \dots < t_{N_t} = 1 \quad \text{with} \quad |t_{l+1} - t_l| = h_t,$$

$$0 = s_0 < s_1 < s_2 < \dots, \quad \text{with} \quad |s_{k+1} - s_k| = h_s.$$

How can we observe the convergence as $N \rightarrow \infty$?

From now on, we use the measure notation ρ_N defined by

$$\rho_N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i(s)}.$$

In the best way, one can consider, for some $\rho \in \mathcal{P}(\Omega_\theta)$ and some metric d ,

$$d(\rho_N, \rho) \lesssim N^{-\alpha}, \quad \text{for some } \alpha > 0.$$

In numeric, it produces **high-cost** in computation. Instead, we would like to see the convergence indirectly using cost function E . Roughly speaking, we want to observe,

$$|E(\rho) - E(\rho_N)| \lesssim N^{-\beta}, \quad \text{for some } \beta > 0.$$

Which β ?

Lemma

If f is bounded and Lipschitz, then we have

$$|E(\rho) - E(\rho_N)| \lesssim \frac{1}{\sqrt{N}}.$$

Proof. We use

$$\|\rho - \bar{\rho}_N\|_{w_1} \lesssim \frac{1}{\sqrt{N}} \quad \text{and} \quad \|\bar{\rho}_N - \rho_N\|_{w_1} \lesssim \frac{1}{\sqrt{N}}.$$

and Grönwall's inequality to derive

$$|u_\rho(t=1, x) - u_{\rho_N}(t=1, x)| \lesssim \frac{1}{\sqrt{N}}.$$

Thanks to finite spatial measure ($|\Omega_x| < \infty$), we get the desired result. \square

Numerical Results

What do we observe in simulation?

Goal: Observe $|E(\rho) - E(\rho_N)| \lesssim \frac{1}{\sqrt{N}}$.

We consider a scalar

$$\alpha_N(s) := \frac{\sqrt{2N}|E(\rho_N) - E(\rho_{2N})|}{\sqrt{2} - 1} \quad (6)$$

and we claim that $\alpha < C$ for any N .

If it is true, then

$$|E(\rho_N) - E(\rho)| < \frac{C}{\sqrt{N}}.$$

We want to show in simulation that $\alpha_N(s) < C$ for all N . For example,

$$\alpha_{25}, \alpha_{50}, \alpha_{100}, \alpha_{200} < C.$$

Test Cases and Setup

- Goal: train $\{\theta_i\}_{i=1}^N$ to recover the true distribution $\rho_*(\theta) = \mathcal{N}(0, 1)$.
- We approximate ρ_* with $\{\theta_i^*\}_{i=1}^{N^*}$, where $N^* = 3000$ and $\theta_i^* \sim \rho_*$.
- We generate the target data $\{(x_j, u_j)\}_{j=1}^{M^*}$ with $M^* = 1000$, where $u_j = u_*(t = 1, x_j)$.
- We then randomly select $M = 500$ samples from the target data set $\{(x_j, u_j)\}_{j=1}^{M^*}$ to train on.
- We initialize $\{\theta_i\}_{i=1}^N$ such that $\theta_i \sim \mathcal{U}[0, 1]$.
- For $f(u, \theta)$, we choose radial basis functions (RBF). That is $f(u, \theta) = \phi(|u - \theta|)$. Specifically we test $f(u, \theta) = \exp(-|u - \theta|^2)$ and $f(u, \theta) = \frac{1}{1 + |u - \theta|^2}$.
- We use a forward Euler timestep of $h_t = 0.1$. We advance Θ_N by its flow map using a step size of $h_s = 100$ (large, but results show it is not unstable). We iterate the system 10,000 times to allow Θ_N to stabilize for all $N = 25, 50, 100, 200, 400$.
- For each N we run thirty trials and average the results to obtain the final parameter set Θ_N .

Test Case: $f(u, \theta) = \exp(-(u - \theta)^2)$

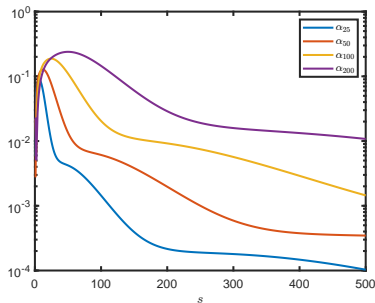
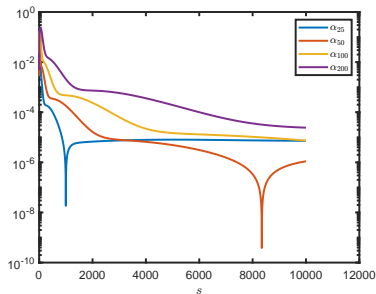


Figure: α (y-axis) versus pseudo-time s (x-axis) for $N = 25, 50, 100, 200$. The right plot zooms in on the region $s \in [0, 500]$.

Test Case: $f(u, \theta) = \exp(-(u - \theta)^2)$

$N = 400$

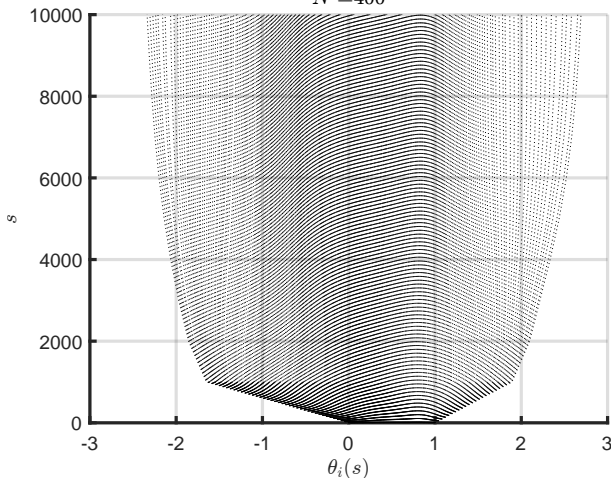
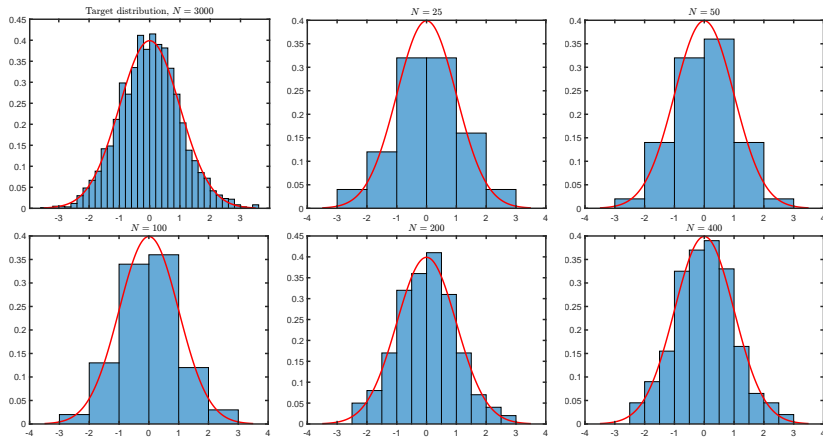


Figure: Trajectories of θ_i (x-axis) as a function of pseudo-time s (y-axis), $N = 400$, $\theta(s = 0) \sim \mathcal{U}[0, 1]$, $\rho_* = \mathcal{N}(0, 1)$.

Test Case: $f(u, \theta) = \exp(-(u - \theta)^2)$



$$\text{Test Case: } f(u, \theta) = \frac{1}{1+(u-\theta)^2}$$

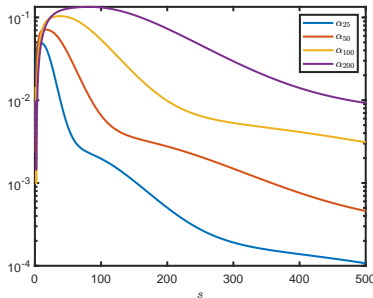
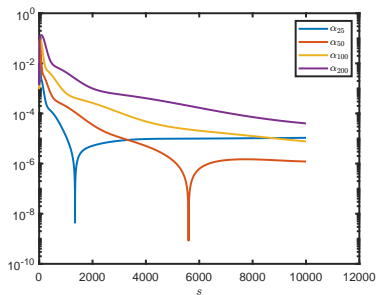


Figure: α (y-axis) versus pseudo-time s (x-axis) for $N = 25, 50, 100, 200$. The right plot zooms in on the region $s \in [0, 500]$.

Test Case: $f(u, \theta) = \frac{1}{1+(u-\theta)^2}$

$N = 400$

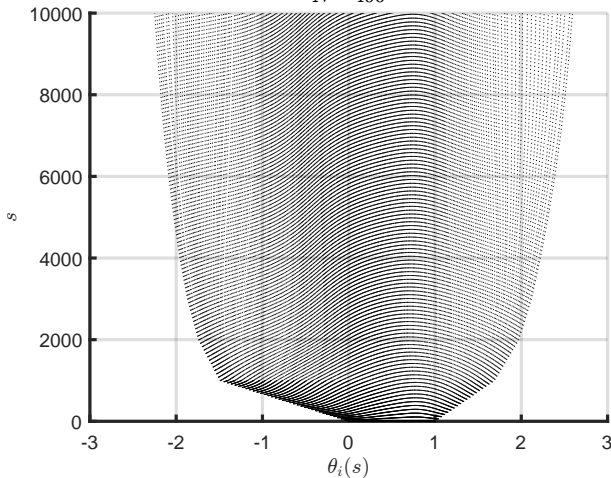


Figure: Trajectories of θ_i (x-axis) as a function of pseudo-time s (y-axis), $N = 400$, $\theta(s=0) \sim \mathcal{U}[0, 1]$, $\rho_* = \mathcal{N}(0, 1)$.

$$\text{Test Case: } f(u, \theta) = \frac{1}{1+(u-\theta)^2}$$

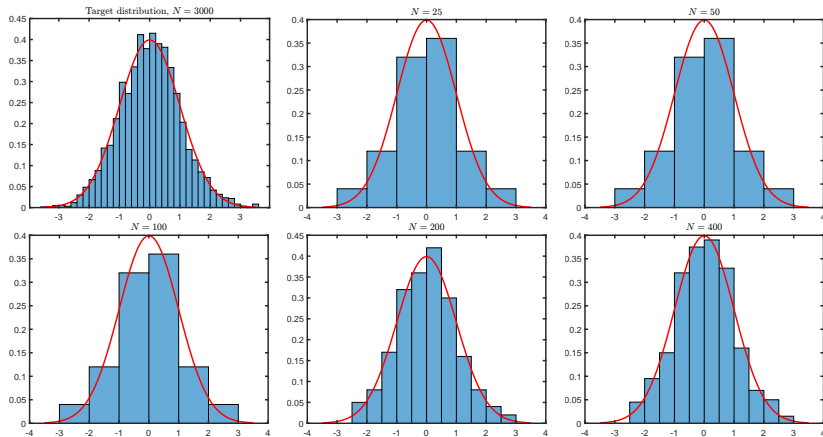


Figure: Convergence $\rho_N \rightarrow \rho_* = \mathcal{N}(0, 1)$, $N = 25, 50, 100, 200, 400$.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 28/28