

A APPENDIX

We add the proof of Theorem 1 and additional numerical experiments here.

A.1 PRELIMINARIES FROM OPTIMAL TRANSPORT THEORY

Definition 2. Suppose X is a metric space equipped with the metric $d(\mathbf{x}, \mathbf{y})$, and μ and ν are two probability measures on X . The *Wasserstein distance* (as known as the *Kantorovich–Rubinstein metric*) $d_{W^d}(\mu, \nu)$ between two probability measures μ and ν for the metric function $d(\mathbf{x}, \mathbf{y})$ is defined to be

$$d_{W^d}(\mu, \nu) = \inf_{\pi \in \Pi(X \times X)} \int_{X \times X} d(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}),$$

where $\Pi(X \times X)$ is the collection of all probability measure on $X \times X$ such that

$$\pi(A \times X) = \mu(A), \quad \pi(X \times B) = \nu(B)$$

for all measurable sets $A, B \subset X$.

For the analysis of the adaptive algorithm in this work, we consider the metric $d_M(\mathbf{x}, \mathbf{y})$ induced by the Euclidean metric $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$

$$d_M(\mathbf{x}, \mathbf{y}) = \min\{M, d(\mathbf{x}, \mathbf{y})\}, \quad \mathbf{x}, \mathbf{y} \in X.$$

Then the metric $d_M(\mathbf{x}, \mathbf{y})$ is always bounded by M (reachable, namely $\|d_M\|_\infty = M$). We denote the Wasserstein distance for $d_M(\mathbf{x}, \mathbf{y})$ by $d_{W^M}(\cdot, \cdot)$.

According to the optimal transport theory, the Wasserstein distance can be described by its dual form (see e.g. Villani (2003), Theorem 1.14 and Remark 1.15 on Page 34).

Theorem 3 (Kantorovich–Rubinstein theorem). *Let X be a Polish space and let d be a lower semi-continuous metric on X . Let $\|\cdot\|_{Lip}$ denote the Lipschitz norm of a function defined as*

$$\|\phi\|_{Lip} = \sup_{\mathbf{x} \neq \mathbf{y}} \frac{|\phi(\mathbf{x}) - \phi(\mathbf{y})|}{d(\mathbf{x}, \mathbf{y})}.$$

Then

$$d_{W^M}(\mu, \nu) = \sup \left\{ \int_X \phi(\mathbf{x}) d(\mu - \nu)(\mathbf{x}) \mid 0 \leq \phi(\mathbf{x}) \leq \|d_M\|_\infty = M, \text{ and } \|\phi\|_{Lip} \leq 1 \right\}.$$

In this work, we restrict ourselves on a compact domain $X = \Omega \subset \mathbb{R}^D$ of learning, and without loss of generality, we assume the Lebesgue measure of Ω is 1.

A.2 THE FIRST CONVERGENCE THEOREM AND ITS PROOF

Theorem 4. *Let μ be the Lebesgue measure on X , which represents the uniform probability distribution on Ω . In addition, we assume Assumption A1 holds.*

Then the optimal value of the min-max problem equation 5 is 0. Moreover, there is a sequence $\{u_n\}_{n=1}^\infty$ of functions with $r(u_n) \neq 0$ for all n , such that it is an optimization sequence of problem equation 5, namely,

$$\lim_{n \rightarrow \infty} \mathcal{J}(u_n, p_n) = 0. \quad (15)$$

for some sequence of functions $\{p_n\}_{n=1}^\infty \subset V$. Meanwhile, this optimization sequence has the following two properties:

1. *The residual sequence $\{r(u_n)\}_{n=1}^\infty$ of $\{u_n\}_{n=1}^\infty$ converges to 0 in $L^2(d\mu)$.*
2. *The renormalized squared residual distributions*

$$d\nu_n \triangleq \frac{r^2(u_n)}{\int_\Omega r^2(u_n(\mathbf{x})) d\mathbf{x}} d\mu \quad (16)$$

converge to the uniform distribution μ in the Wasserstein distance d_{W^M} .

Proof. Consider a minimizing sequence $u_n, n = 1, 2, \dots$ of

$$\inf_u \int_{\Omega} r^2(u(\mathbf{x})) d\mathbf{x}, \quad (17)$$

where without loss of generality, we can assume that $\int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x} \leq \frac{1}{n}$.

Now

$$\begin{aligned} & \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \mathcal{J}(u_n, p) \\ &= \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \left[\int_{\Omega} r^2(u_n(\mathbf{x})) p(\mathbf{x}) d\mathbf{x} - \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x} \int_{\Omega} p(\mathbf{x}) d\mathbf{x} + \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x} \int_{\Omega} p(\mathbf{x}) d\mathbf{x} \right] \\ &\leq \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x} \left(\sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \left[\int_{\Omega} p(\mathbf{x}) d\nu_n(\mathbf{x}) - \int_{\Omega} p(\mathbf{x}) d\mathbf{x} \right] + \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \int_{\Omega} p(\mathbf{x}) d\mathbf{x} \right) \\ &= (d_{WM}(\nu_n, \mu) + M) \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x}. \end{aligned} \quad (18)$$

By the assumption of the theorem, for each n , we can find a function $\tilde{u}_n(\mathbf{x})$ so that the Wasserstein distance $d_{WM}(\tilde{\nu}_n, \mu) \leq \frac{1}{n}$, where $\tilde{\nu}_n$ is the measure defined as in equation 16 by replacing $u_n(\mathbf{x})$ with $\tilde{u}_n(\mathbf{x})$. In fact, for each n , we can find, by partition of unity, a sequence of functions in $C_c^\infty(\Omega)$ converging to $\mathbb{1}_{\Omega}$ in the Sobolev norm of $W^{k,1}$ (See for example Evans (2010)). So we can find a function w_n in $C_c^\infty(\Omega)$, such that $\|w_n(\mathbf{x}) - \mathbb{1}_{\Omega}(\mathbf{x})\|_1 \leq \frac{1}{n}$ on Ω . Since r is a surjection, there is some $\tilde{u}_n(\mathbf{x})$ so that

$$r^2(\tilde{u}_n) = w_n \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x},$$

and

$$\begin{aligned} \int_{\Omega} r^2(\tilde{u}_n) d\mathbf{x} &= \int_{\Omega} w_n(\mathbf{x}) d\mathbf{x} \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x} \\ &\leq (1 + \int_{\Omega} \mathbb{1}_{\Omega}(\mathbf{x}) d\mathbf{x}) \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x} \\ &= 2 \int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x}. \end{aligned}$$

This means $\{\tilde{u}_n\}_{n=1}^\infty$ is also a minimizing sequence of equation 17, and it yields

$$\begin{aligned} d_{WM}(\tilde{\nu}_n, \mu) &= \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \left[\int_{\Omega} p(\mathbf{x}) d\tilde{\nu}_n(\mathbf{x}) - \int_{\Omega} p(\mathbf{x}) d\mathbf{x} \right] \\ &= \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \int_{\Omega} p(\mathbf{x}) \left[\frac{r^2(\tilde{u}_n)(\mathbf{x})}{\int_{\Omega} r^2(\tilde{u}_n(\mathbf{x})) d\mathbf{x}} - \mathbb{1}_{\Omega}(\mathbf{x}) \right] d\mathbf{x} \\ &= \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \int_{\Omega} p(\mathbf{x}) [w_n(\mathbf{x}) - \mathbb{1}_{\Omega}(\mathbf{x})] d\mathbf{x} \\ &\leq \frac{M}{n}. \end{aligned}$$

So we get from equation 18 that

$$0 \leq \lim_{n \rightarrow \infty} \sup_{\substack{\|p\|_{\text{Lip}} \leq 1 \\ 0 \leq p \leq M}} \mathcal{J}(\tilde{u}_n, p) \leq \lim_{n \rightarrow \infty} 4M \int_{\Omega} r^2(u_n) d\mathbf{x} = 0,$$

which means that $\{\tilde{u}_n\}_{n=1}^\infty$ is also a minimizing sequence of equation 5, that is,

$$\lim_{n \rightarrow \infty} \mathcal{J}(\tilde{u}_n, p_n) = 0, \quad (15)$$

for some sequence of functions $\{p_n\}_{n=1}^\infty \subset V$. Meanwhile, we have the following properties of \tilde{u}_n :

1. The residual sequence $\{r(\tilde{u}_n)\}_{n=1}^{\infty}$ converges to 0 in $L^2(d\mu)$, since

$$\int_{\Omega} r^2(\tilde{u}_n) d\mathbf{x} \leq 2 \int_{\Omega} r^2(u_n) d\mathbf{x} \leq \frac{2}{n} \rightarrow 0, \quad \text{as } n \rightarrow \infty$$

2. The renormalized squared residual distributions

$$d\tilde{\nu}_n \triangleq \frac{r^2(\tilde{u}_n)}{\int_{\Omega} r^2(\tilde{u}_n(\mathbf{x})) d\mathbf{x}} d\mu$$

converges to the uniform distribution μ in the Wasserstein distance d_{WM} .

□

A.3 REPLACEMENT OF THE BOUNDEDNESS CONDITION IN THEOREM 4

For the boundedness constraint for “test function” p in 4, we prove that it can be removed in our circumstance. And with the following lemma and its following remark, and Theorem 4, we can obtain our main Theorem 1, which is stated again with its assumption in the following.

Assumption. The operator r in equation 7 is a surjection from a function space $E_1(\mathbb{R}^D)$ to $C_c^\infty(\Omega)$, the class of C^∞ functions that are compactly supported on Ω .

Theorem. Let μ be the Lebesgue measure on \mathbb{R}^D , which represents the uniform probability distribution on Ω . In addition, we assume Assumption A1 holds. Then the optimal value of the min-max problem equation 7 is 0. Moreover, there is a sequence $\{u_n\}_{n=1}^{\infty}$ of functions with $r(u_n) \neq 0$ for all n , such that it is an optimization sequence of equation 7, namely,

$$\lim_{n \rightarrow \infty} \mathcal{J}(u_n, p_n) = 0,$$

for some sequence of functions $\{p_n\}_{n=1}^{\infty}$ satisfying the constraints in equation 7. Meanwhile, this optimization sequence has the following two properties:

1. The residual sequence $\{r(u_n)\}_{n=1}^{\infty}$ of $\{u_n\}_{n=1}^{\infty}$ converges to 0 in $L^2(d\mu)$.
2. The renormalized squared residual distributions

$$d\nu_n \triangleq \frac{r^2(u_n)}{\int_{\Omega} r^2(u_n(\mathbf{x})) d\mathbf{x}} d\mu(\mathbf{x})$$

converge to the uniform distribution μ in the Wasserstein distance d_{WM} .

Although the residue r^2 is renormalized to a probability distribution for the analysis of the algorithm, itself is not a probability distribution, and not treated as so. Actually, in the implementation of our algorithm, the “test function” p is seen as sampling distribution density and the residue r^2 is just the PDE operator (or any kind of objective function whose minimum is 0). In the implementation, we establish p as a generative model, that is, an invertible transform between an unknown distribution (an adversarial distribution to the residual distribution if we think the algorithm as a similarity to GANs) and an “easy-to-sample” distribution such as normal or uniform distribution. So we assume p to be the density function of a probability distribution. Under this assumption, we have the following result.

Lemma 5. Let Ω be a compact subset of \mathbb{R}^D . If a positive function $f : \Omega \rightarrow \mathbb{R}$ is K -Lipschitz continuous, and f is the density function of a probability distribution, namely, $\int_{\Omega} f d\mathbf{x} = 1$, then there is some constant $M = M(\Omega, K)$, so that $f \leq M$. In other words,

$$f \leq M, \quad \forall f \in \mathcal{S} = \{f \geq 0 \mid \|f\|_{Lip} \leq K, \text{ and } \int_{\Omega} f d\mathbf{x} = 1\}.$$

Proof. For any $x, y \in \Omega$, we have

$$0 \leq f(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{y}) + f(\mathbf{y}) \leq K|\mathbf{x} - \mathbf{y}| + f(\mathbf{y}) \leq K\mathcal{D}(\Omega) + f(\mathbf{y}),$$

where $\mathcal{D}(\Omega)$ is the diameter of Ω . Taking integral with respect to \mathbf{y} over Ω on both sides, we have

$$0 \leq f(\mathbf{x})\mu(\Omega) \leq K\mathcal{D}(\Omega)\mu(\Omega) + 1,$$

where $\mu(\Omega)$ is the Lebesgue measure (volume) of Ω , that is,

$$0 \leq f(\mathbf{x}) \leq K\mathcal{D}(\Omega) + \frac{1}{\mu(\Omega)}.$$

So we have

$$M = M(\Omega, K) = K\mathcal{D}(\Omega) + \frac{1}{\mu(\Omega)}.$$

□

e The converse of this lemma is also true in the sense that if f is bounded by some constant M , then the integral $\int_{\Omega} f d\mathbf{x} \leq M\mu(\Omega)$, and f can be renormalized into a probability density function with constant $M\mu(\Omega)$. And similar to boundedness for the gradient (or Lipschitz constant) discussed in section 3.3, a constant renormalizer will not affect the training procedure.

A.4 DEVIATION OF EQUATION EQUATION 9 AND IT SOLUTION

For a given $r(\mathbf{x}; \boldsymbol{\theta})$, consider the following minimization problem:

$$\min_{p_{\alpha} > 0} \mathcal{L}(p_{\alpha}) = \beta \int_{\Omega} |\nabla_{\mathbf{x}} p_{\alpha}|^2 d\mathbf{x} - \int_{\Omega} r^2(\mathbf{x}; \boldsymbol{\theta}) p_{\alpha}(\mathbf{x}) d\mathbf{x} + \lambda \left(\int_{\Omega} p_{\alpha}(d\mathbf{x}) - 1 \right),$$

where the positivity of p_{α} is guaranteed by the KRnet and λ is the Lagrange multiplier for the mass conservation of PDF. Assuming that $\frac{\partial p_{\alpha}}{\partial \mathbf{n}} = 0$ on the boundary $\partial\Omega$, where \mathbf{n} is a unit normal vector on $\partial\Omega$ pointing outward. We have the first-order variation of $\mathcal{L}(p_{\alpha})$ for a perturbation function $\delta p(\mathbf{x})$

$$\begin{aligned} \delta \mathcal{L} &= 2\beta \int_{\Omega} \nabla p_{\alpha} \cdot \nabla \delta p d\mathbf{x} - \int_{\Omega} r^2 \delta p d\mathbf{x} + \lambda \int_{\Omega} \delta p d\mathbf{x} \\ &= 2\beta \left(\int_{\partial\Omega} \delta p \nabla p_{\alpha} \cdot \mathbf{n} d\Gamma - \int_{\Omega} \delta p \nabla^2 p_{\alpha} d\mathbf{x} \right) - \int_{\Omega} r^2 \delta p d\mathbf{x} + \lambda \int_{\Omega} \delta p(\mathbf{x}) d\mathbf{x} \\ &= -2\beta \int_{\Omega} \delta p \nabla^2 p_{\alpha} d\mathbf{x} - \int_{\Omega} r^2 \delta p d\mathbf{x} + \lambda \int_{\Omega} \delta p(\mathbf{x}) d\mathbf{x} \\ &= - \int_{\Omega} (2\beta \nabla^2 p_{\alpha} + r^2 - \lambda) \delta p d\mathbf{x}, \end{aligned}$$

where we applied integration by parts and the homogeneous Neuman boundary conditions. The optimality condition $\frac{\delta \mathcal{L}}{\delta p} = 0$ yields

$$\begin{cases} 2\beta \nabla^2 p_{\alpha}(\mathbf{x}) + r^2(\mathbf{x}; \boldsymbol{\theta}) - \lambda = 0, & \mathbf{x} \in \Omega, \\ \frac{\partial p_{\alpha}}{\partial \mathbf{n}} = 0, & \mathbf{x} \in \partial\Omega. \end{cases} \quad (19)$$

From the compatibility condition for Neumann problems, we have

$$\int_{\Omega} (r^2(\mathbf{x}; \boldsymbol{\theta}) - \lambda) d\mathbf{x} = 0, \quad (20)$$

which yields that

$$\lambda = \frac{1}{|\Omega|} \int_{\Omega} r^2(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}.$$

Assume that Ω is a bounded domain with smooth boundary. It can be shown that if $r \in H^k(\Omega)$ and $\partial\Omega \in C^{k+2}$ with $k \in \mathbb{N}$, the solution of equation equation 9 satisfies Taylor (2011)

$$\|p_{\alpha}\|_{H^{k+2}(\Omega)} \leq C \|f\|_{H^k(\Omega)},$$

where $f(\mathbf{x}) = (\lambda - r^2)/(2\beta)$ and $C > 0$ is a general constant that does not depend on r . According to the Sobolev Imbedding Theorem Adams & John Fournier (2003),

$$W^{k,1}(\Omega) \rightarrow C^{0,1}(\bar{\Omega}),$$

when $D = k - 1$. Thus up to a set of measure zero, we have

$$\|p_\alpha\|_{C^{0,1}(\bar{\Omega})} \leq C_1 \|p_\alpha\|_{W^{k,1}(\Omega)} \leq C_2 \|p_\alpha\|_{H^k(\Omega)},$$

where C_1 and C_2 are general constants independent of p_α . So p_α is Lipschitz continuous when the boundary and $r(\mathbf{x})$ are sufficiently smooth. However, this also means that the H_1 regularization used in equation 8 induces a weaker constraint than the Lipschitz condition in Lemma 5.

A.5 SUPPLEMENTARY EXPERIMENTS

About the setting of $s(\mathbf{x})$ and $g(\mathbf{x})$. The source term $s(\mathbf{x})$ is derived by the exact solution, i.e., we can set the source function by plugging the exact solution into the equation to get $s(\mathbf{x})$. We set $g(\mathbf{x}) = u(\mathbf{x})$ since the Dirichlet boundary condition is imposed on $\partial\Omega$.

Parametric Burgers' Equation. We also test the proposed AAS method using parametric PDEs that are commonly used in the design of engineering systems and uncertainty quantification. Specifically, we consider the following parametric Burgers' equation, which is a benchmark problem studied in DeepXDE.

$$\begin{aligned} \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} &= \nu \left[\left(\frac{\partial u}{\partial x} \right)^2 + \left(\frac{\partial u}{\partial y} \right)^2 \right] \\ \frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} &= \nu \left[\left(\frac{\partial v}{\partial x} \right)^2 + \left(\frac{\partial v}{\partial y} \right)^2 \right] \\ x, y &\in [0, 1], \text{ and } t \in [0, 1] \end{aligned}$$

where u and v are the velocities along x and y directions respectively, and $\nu \in (0, 1]$ is a parameter that represents the kinematic viscosity of fluid. Here, the Dirichlet boundary conditions are imposed on all boundaries. The exact solution is obtained as follows.

$$\begin{aligned} u(x, y, t) &= \frac{3}{4} - \frac{1}{4[1 + \exp((-4x + 4y - t)/(32\nu))]}, \\ v(x, y, t) &= \frac{3}{4} + \frac{1}{4[1 + \exp((-4x + 4y - t)/(32\nu))]}, \end{aligned}$$

The problem setup space is $\mathbf{x} = [t, x, y, \nu]$, i.e., $D = 4$. When ν is small, solving this problem is quite challenging. We use the proposed AAS method to train a neural network $u_\theta(\mathbf{x})$ to approximate the solution over the entire space $\mathbf{x} = [t, x, y, \nu] \in [0, 1]^4$. Figure 6 shows the numerical results, which demonstrate that the proposed AAS method is able to accurately solve this parametric Burgers' equation. We can train the models using the strategy as discussed in Remark 2, i.e., we gradually add the data points to the current training set. AAS with fixed $\beta = 5$ means that we use a similar training strategy as DAS-G presented in Tang et al. (2023) with a fixed β , while AAS with decay $\beta = 5$ means that β has a decay scheme at every 100 stages with decay rate 0.9. Adding the data points gradually to the current set of random samples is more stable than that of replacing all data points.

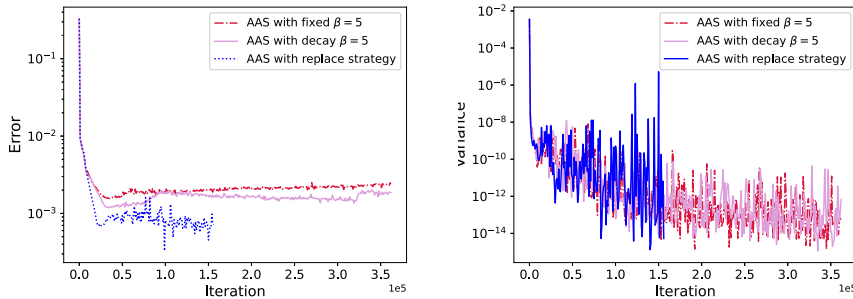


Figure 6: The results of the parametric Burgers' equation. Left: The error behavior. Right: The evolution of the variance.