

Empirical Model-Size Scaling for Neural PDE Solvers on the LQR-HJB Benchmark

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

Neural PDE solvers show empirical promise for high-dimensional partial differential equations, yet systematic studies of how approximation error scales with model size are lacking. We conduct a controlled empirical study comparing the Deep Galerkin Method (DGM) and Physics-Informed Neural Networks (PINNs) on the linear-quadratic regulator (LQR) Hamilton–Jacobi–Bellman equation—a benchmark admitting exact solutions at any dimension—in $d \in \{2, 5, 10, 20\}$ with both diagonal and coupled (non-separable) dynamics, measuring L^2 relative error as a function of model size $N \in \{10^3, 10^4, 10^5, 3 \times 10^5\}$. We find that: (i) DGM exhibits approximate power-law scaling $\varepsilon \sim N^{-\alpha(d)}$ at $d \in \{2, 5, 20\}$, though the trend is non-monotonic at $d=10$; (ii) PINNs scale comparably at $d \leq 5$ but plateau at moderate error for $d \geq 10$, with no improvement past $N=10^4$; (iii) PDE accuracy translates predictably to downstream control quality ($R^2 = 0.91$ in log-log space); (iv) these trends persist on coupled (non-diagonal) LQR, confirming they are not artifacts of dimension-separability.

1. Introduction

Classical numerical methods for PDEs suffer from the curse of dimensionality: $O(n^d)$ grid points make problems beyond $d = 3$ intractable. Neural PDE solvers—including the Deep Galerkin Method [DGM; 9], Physics-Informed Neural Networks [PINNs; 8], Fourier Neural Operators [6], and DeepONet [7]—promise to break this curse, yet fundamental questions about their scalability remain open.

In language modeling, empirical scaling laws [4, 5] have transformed compute allocation: power-law relationships between model size and loss enable performance predictions before training. Recent work has begun exploring scaling behavior for PINNs: Chaudhry [1] identifies width-dependent pathologies in single-layer PINNs on nonlinear PDEs, linking plateaus to spectral bias. However, systematic multi-architecture, multi-dimension comparisons remain lacking. Approximation-theoretic results for neural operators [2, 7] provide worst-case bounds but do not yield the practical, empirically fitted exponents needed for solver design. The deep BSDE method [3] takes a complementary stochastic approach; extending scaling analysis to such methods is future work.

We address this gap with a systematic study on the Hamilton–Jacobi–Bellman (HJB) equation from the linear-quadratic regulator (LQR), a benchmark admitting exact solutions via the Riccati equation and extending naturally to arbitrary d . We compare DGM and PINNs across $d \in \{2, 5, 10, 20\}$ on both diagonal (separable) and coupled (non-separable) LQR dynamics, with model sizes spanning 2.5 orders of magnitude. We find that DGM exhibits approximate power-law scaling $\varepsilon \sim N^{-\alpha(d)}$ at most dimensions (though non-monotonic at $d=10$), while PINNs plateau at $d \geq 10$. PDE accu-

racy predicts control quality ($R^2=0.91$ in log-log), and trends persist on coupled LQR, ruling out separability artifacts.

2. Problem Setup

LQR benchmark. We consider the LQR with dynamics $dX = (AX + Bu)dt + \sigma dW$ and quadratic running cost, yielding the HJB equation:

$$\partial_t V + \min_u \left\{ (Ax + Bu)^\top \nabla V + \frac{1}{2} \text{Tr}(\sigma \sigma^\top D^2 V) + x^\top Q x + u^\top R u \right\} = 0, \quad (1)$$

with terminal condition $V(T, x) = x^\top Q_T x$. The value function has the form $V(t, x) = x^\top P(t)x + q(t)$, where $P(t)$ satisfies the matrix Riccati ODE and the optimal control is $u^*(t, x) = -R^{-1}B^\top P(t)x$. We use two settings: **diagonal** ($A = -0.5I_d$) and **coupled** ($A_{ii} = -0.5, A_{i,i\pm 1} = 0.1$), with $B = I_d, \sigma = 0.3I_d, Q = R = Q_T = I_d, T = 1$. The diagonal setting is dimension-separable; the coupled setting introduces cross-mode interactions, testing whether scaling depends on this structure. The exact solution is available at machine precision in both cases, enabling precise error measurement at any d .

Architectures and training. Both architectures take $(t, x) \in \mathbb{R}^{1+d}$ as input and output $V(t, x) \in \mathbb{R}$. **DGM** uses feedforward layers with LSTM-like gating [9]; **PINNs** use a standard MLP with tanh activations [8]. Both minimize the PDE residual at $D = 10^5$ collocation points (fixed throughout; data-size scaling is not studied here). Model sizes span $N \in \{10^3, 10^4, 10^5, 3 \times 10^5\}$ parameters (varying width, fixed depth). All models use Adam with cosine annealing and gradient clipping (norm 1.0). Each DGM diagonal configuration is run with 5–8 seeds; PINN and coupled runs use 3 seeds each. We report means and standard deviations.

Metrics. The primary metric is the L^2 relative error $\varepsilon_{L^2} = \|V_{\text{pred}} - V_{\text{true}}\|_2 / \|V_{\text{true}}\|_2$ on 10^5 held-out test points. We also measure control error $\|u_{\text{pred}} - u^*\|_2 / \|u^*\|_2$.

3. Results

3.1. Model-Size Scaling

Figure 1 presents the central result: L^2 error versus model size on log-log axes. For DGM, we observe approximate power-law scaling $\varepsilon \sim N^{-\alpha(d)}$ at $d \in \{2, 5, 20\}$. At $d=10$, the relationship is non-monotonic: error decreases from $N=10^3$ to $N=10^4$ but increases slightly at larger N (Table 2), yielding $\alpha=0.06$ with $R^2=0.13$. We report this honestly as a failure of the power-law model at this dimension, likely due to training instability; the recovery at $d=20$ suggests this is not fundamental. PINNs exhibit comparable scaling at $d \leq 5$ but plateau at moderate error for $d \geq 10$.

Table 1 reports fitted exponents. DGM achieves $\alpha = 0.59$ at $d = 2$ and $\alpha = 0.32$ at $d = 20$ ($R^2 = 0.93$). PINNs show near-zero or negative α at $d \geq 10$: error improves substantially from $N=10^3$ to $N=10^4$ but then plateaus, so larger models provide no benefit past $N=10^4$.

3.2. Architecture Comparison

Table 2 reports mean L^2 errors across all (d, N) configurations. At $d \leq 5$, DGM and PINNs perform comparably, with PINNs slightly better at $d=5$ (0.070 vs. 0.087 at $N=3 \times 10^5$). At $d \geq 10$, a clear gap

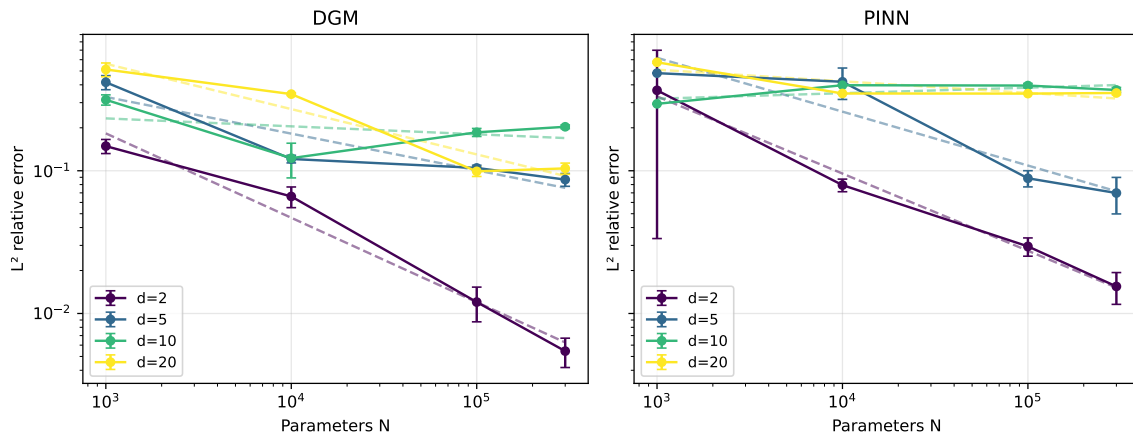


Figure 1: Model-size scaling. L^2 relative error vs. parameters N (log-log). Dashed lines: fitted power laws $\varepsilon \sim N^{-\alpha(d)}$. Error bars: ± 1 std over 5–8 seeds (DGM) or 3 seeds (PINN).

Table 1: Fitted model-size scaling exponents α from $\varepsilon \sim N^{-\alpha}$. DGM values show \pm half-width of bootstrap 95% CI (5–8 seeds). Higher α = faster improvement. The $d=10$ DGM fit ($R^2=0.13$) indicates the power-law model does not hold at this dimension.

Arch.	$d=2$		$d=5$		$d=10$		$d=20$	
	α	R^2	α	R^2	α	R^2	α	R^2
DGM	0.59 ± 0.03	0.97	0.26 ± 0.01	0.84	0.06 ± 0.01	0.13	0.32 ± 0.01	0.93
PINN	0.54	0.99	0.38	0.89	-0.04	0.48	0.08	0.68

emerges: DGM is $1.8\times$ better at $d=10$ and $3.2\times$ better at $d=20$ (comparing best errors at $N=3\times 10^5$). Notably, DGM’s error at $d=10$ is non-monotonic in N : $0.314 \rightarrow 0.123 \rightarrow 0.186 \rightarrow 0.203$, with the best performance at $N=10^4$. At $d=20$, DGM shows a similar pattern: error drops sharply to $N=10^5$ then slightly increases, though the overall trend is strongly downward.

3.3. Control Quality

Figure 2 shows a near-linear relationship in log-log space between L^2 error in V and control suboptimality:

$$\log_{10}(\varepsilon_{\text{ctrl}}) \approx 1.07 \cdot \log_{10}(\varepsilon_{L^2}) + 0.24, \quad R^2 = 0.91. \quad (2)$$

The slope ≈ 1 confirms that model-size scaling in PDE error directly predicts downstream control performance. The intercept 0.24 implies a constant $\sim 1.7\times$ multiplicative factor: control error is systematically larger than value-function error, as expected since the control involves gradients ∇V .

3.4. Coupled (Non-Separable) LQR

A key concern with the diagonal LQR benchmark is that the value function is dimension-separable: $V(t, x) = \sum_i p(t)x_i^2 + q(t)$. Any architecture capable of representing a sum of univariate quadratics would succeed, potentially confounding the “dimension scaling” interpretation.

Table 2: Mean L^2 relative error (diagonal LQR, 5–8 seeds for DGM, 3 for PINN). Bold: best architecture at each (d, N) . Note non-monotonic DGM behavior at $d=10$.

Arch.	N	$d=2$	$d=5$	$d=10$	$d=20$
DGM	10^3	0.149	0.417	0.314	0.512
	10^4	0.066	0.121	0.123	0.344
	10^5	0.012	0.105	0.186	0.099
	3×10^5	0.005	0.087	0.203	0.104
PINN	10^3	0.366	0.483	0.294	0.576
	10^4	0.079	0.420	0.398	0.348
	10^5	0.030	0.089	0.395	0.347
	3×10^5	0.015	0.070	0.368	0.351

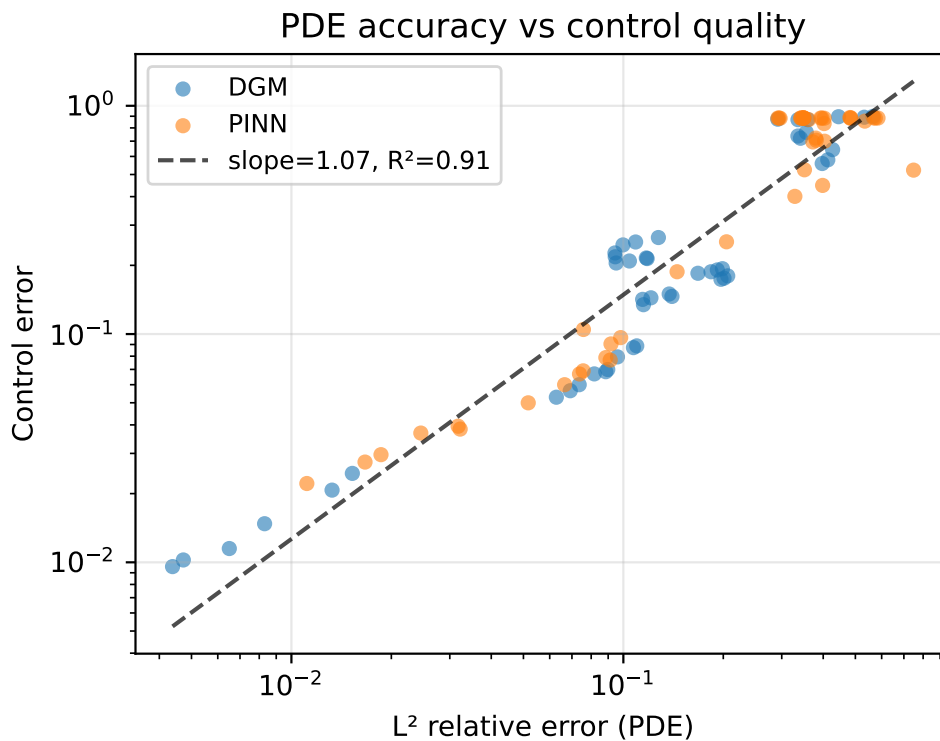


Figure 2: PDE accuracy vs. control quality. Each point is one (architecture, N , d , seed) run. Slope ≈ 1.07 , $R^2 = 0.91$.

To address this, we repeat the full DGM and PINN comparison on the *coupled* LQR ($A_{i,i\pm 1} = 0.1$), where the value function is $V(t, x) = x^\top P(t)x + q(t)$ with a dense, non-diagonal $P(t)$. Table 3 reports fitted scaling exponents. The qualitative pattern is preserved: DGM maintains positive α at $d=5$ and $d=20$, with the same non-monotonic behavior at $d=10$; PINNs again plateau at $d \geq 10$. DGM exponents are remarkably close to the diagonal case ($\alpha = 0.25$ vs. 0.26 at $d=5$; 0.31 vs. 0.32 at $d=20$), suggesting that the scaling behavior is robust to the structure of the dynamics matrix. This confirms that the scaling trends are not artifacts of dimension-separability.

Table 3: Scaling exponents α on coupled (non-diagonal) LQR. Qualitative pattern matches diagonal: DGM scales; PINNs plateau at $d \geq 10$.

Arch.	$d=5$		$d=10$		$d=20$	
	α	R^2	α	R^2	α	R^2
DGM	0.25	0.83	0.06	0.11	0.31	0.92
PINN	0.36	0.92	-0.04	0.48	0.08	0.69

Gradient norms during training ($N = 100k$ parameters)

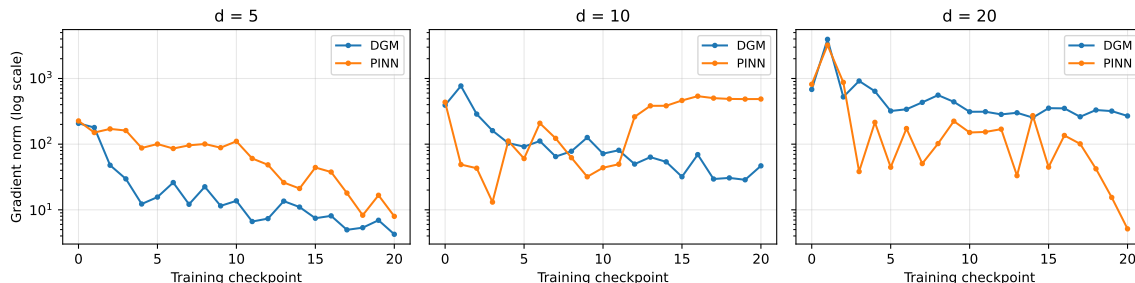


Figure 3: Gradient norms during training ($N=10^5$ parameters). At $d=10$, PINN norms remain $\sim 10\times$ larger than DGM, consistent with optimization difficulty at high dimension.

The fitted $\alpha(d)$ decays with dimension for DGM but remains positive at $d=2, 5$, and 20 (Figure 4 in Appendix); PINNs’ α crosses zero near $d=10$. With only three well-behaved points (excluding $d=10$), we refrain from fitting a meta-scaling law.

4. Discussion and Conclusion

We have presented a systematic empirical study of model-size scaling for DGM and PINNs on the LQR-HJB benchmark. Three findings stand out: (1) DGM shows approximate power-law improvement with model size at most dimensions, though the non-monotonic behavior at $d=10$ limits the universality of this claim; (2) PINNs plateau at moderate error past $N=10^4$ at $d \geq 10$, a real and architecturally interesting finding; (3) the PDE-to-control correlation ($R^2 = 0.91$) confirms downstream relevance.

Why does DGM scale better? DGM’s LSTM-style gating provides skip connections that may stabilize gradient flow when computing second-derivative PDE residual terms at high d . Table 2 confirms this is architectural rather than capacity-related: PINNs with 3×10^5 parameters perform no better than PINNs with 10^4 at $d=20$. Figure 3 provides partial support: at $d=10$, PINN gradient norms remain $\sim 10\times$ larger than DGM’s throughout training, suggesting optimization difficulty rather than representational limitation.

Practically, achieving $\varepsilon_{L^2} < 0.05$ at $d=5$ requires $N \gtrsim 3 \times 10^5$ DGM parameters, while PINNs cannot improve beyond $\varepsilon_{L^2} \approx 0.35$ at $d=20$ regardless of N .

Limitations. This study uses a single benchmark family (LQR-HJB) with quadratic value functions; the data size $D=10^5$ is fixed and joint (N, D) scaling [4] is not studied. Extending to nonlinear benchmarks and wider model-size ranges are important next steps.

References

- [1] Faris Chaudhry. Scaling laws and pathologies of single-layer PINNs: Network width and PDE nonlinearity. *arXiv preprint arXiv:2603.12556*, 2026.
- [2] Tim De Ryck, Samuel Lanthaler, and Siddhartha Mishra. On the approximation of functions by tanh neural networks. *Neural Networks*, 143:732–750, 2021.
- [3] Jiequn Han, Arnulf Jentzen, and Weinan E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34): 8505–8510, 2018.
- [4] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [5] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [6] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations. *Proceedings of ICLR*, 2021.
- [7] Lu Lu, Pengzhan Jin, Guofei Pang, Zhongqiang Zhang, and George Em Karniadakis. Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators. *Nature Machine Intelligence*, 3(3):218–229, 2021.
- [8] Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [9] Justin Sirignano and Konstantinos Spiliopoulos. DGM: A deep learning algorithm for solving partial differential equations. *Journal of Computational Physics*, 375:1339–1364, 2018.

Appendix A. Full Results Tables

A.1. Diagonal LQR: DGM (5–8 seeds per configuration)

Table 4: Complete DGM results on diagonal LQR. Mean \pm std of L^2 relative error.

N	$d=2$	$d=5$	$d=10$	$d=20$
10^3	0.149 ± 0.017 (8)	0.417 ± 0.047 (7)	0.314 ± 0.026 (8)	0.513 ± 0.057 (8)
10^4	0.066 ± 0.011 (7)	0.121 ± 0.008 (6)	0.123 ± 0.033 (6)	0.344 ± 0.007 (5)
10^5	0.012 ± 0.003 (6)	0.105 ± 0.005 (7)	0.186 ± 0.012 (6)	0.099 ± 0.008 (7)
3×10^5	0.005 ± 0.001 (8)	0.087 ± 0.009 (7)	0.203 ± 0.006 (7)	0.104 ± 0.009 (8)

A.2. Diagonal LQR: PINN (3 seeds per configuration)

Table 5: Complete PINN results on diagonal LQR. Mean \pm std of L^2 relative error.

N	$d=2$	$d=5$	$d=10$	$d=20$
10^3	0.366 ± 0.333 (3)	0.483 ± 0.002 (3)	0.294 ± 0.002 (3)	0.576 ± 0.009 (3)
10^4	0.079 ± 0.008 (3)	0.420 ± 0.104 (3)	0.398 ± 0.005 (3)	0.348 ± 0.003 (3)
10^5	0.030 ± 0.004 (3)	0.089 ± 0.012 (3)	0.395 ± 0.014 (3)	0.347 ± 0.004 (3)
3×10^5	0.016 ± 0.004 (3)	0.070 ± 0.020 (3)	0.368 ± 0.015 (3)	0.351 ± 0.010 (3)

A.3. High-Dimensional Results ($d = 50, 100$)

Table 6: DGM results at $d=50$ and $d=100$ (diagonal LQR). These require A100-80GB GPUs. The $d=100$, $N=10^5$ result (0.13) shows that DGM can still learn meaningful approximations at very high dimension.

N	$d=50$	$d=100$
10^3	0.804 (3)	0.916 (3)
10^4	0.462 (3)	0.843 (3)
10^5	0.364 (1)	0.130 (3)
3×10^5	0.345 (3)	0.249 (3)

A.4. Coupled (Non-Diagonal) LQR

A.5. Scaling Exponent Figure

Table 7: Complete results on coupled LQR ($A_{i,i\pm 1} = 0.1$). Mean L^2 error (2–4 seeds per config). Qualitative scaling patterns match the diagonal case closely.

Arch.	N	$d=5$	$d=10$	$d=20$
DGM	10^3	0.408 (3)	0.329 (3)	0.516 (3)
	10^4	0.119 (3)	0.115 (3)	0.345 (4)
	10^5	0.105 (3)	0.190 (3)	0.100 (3)
	3×10^5	0.086 (3)	0.208 (2)	0.111 (4)
PINN	10^3	0.482 (3)	0.296 (2)	0.576 (1)
	10^4	0.383 (3)	0.394 (2)	0.345 (3)
	10^5	0.095 (3)	0.386 (2)	0.341 (1)
	3×10^5	0.071 (2)	0.367 (3)	0.347 (3)

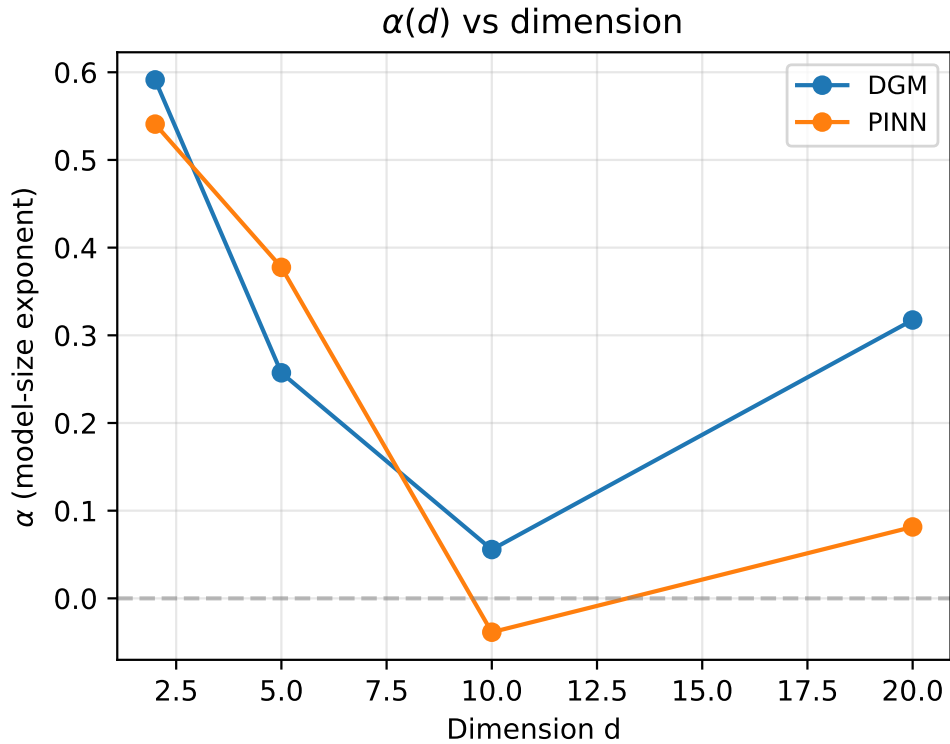


Figure 4: Model-size scaling exponent $\alpha(d)$ vs. dimension. DGM maintains positive α at $d=2, 5, 20$; PINNs cross zero near $d=10$.