
Data Attribution for Model/Learning Based Control (Extended Abstract)

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Ensuring safety of learning-based robotic controllers requires understanding which
2 training data influences performance, yet identifying influential trajectories through
3 exhaustive retraining is computationally prohibitive. We introduce a framework using
4 influence functions to efficiently approximate the impact of individual training
5 trajectories on learned dynamics and control performance. We formulate IF1 to
6 estimate effects on model accuracy and IF2 to quantify impacts on LQR control
7 cost—a proxy for tracking error and stability. Empirical validation demonstrates
8 strong correlations between influence predictions and ground truth.

1 Introduction

9 The deployment of learning-based control systems in safety-critical robotic applications—from
10 autonomous vehicles Levinson et al. [2011] to surgical robots Shademan et al. [2016]—requires
11 rigorous quality assurance of training data. Corrupted or low-quality training trajectories can propa-
12 gate through the learning pipeline, resulting in controllers that exhibit unsafe behaviors: excessive
13 tracking errors, oscillations, or instability Dean et al. [2020]. Understanding which training data most
14 influences controller performance is critical for ensuring reliable deployment Amodei et al. [2016].

15 Traditional LOO retraining for assessing data importance is computationally prohibitive for realistic
16 datasets with thousands of trajectories. Recent advances in machine learning and large language
17 models demonstrate remarkable success using influence functions for efficient data attribution
18 Grosse et al. [2023], including identifying mislabeled examples Koh and Liang [2017], detecting
19 memorization Feldman and Zhang [2020], and guiding data curation Hammoudeh and Lowd [2022].
20 This motivates applying influence functions to learning-based control, where understanding how
21 training trajectories impact controller performance is critical for safety-critical robotic applications.

22 **Our Contribution.** We introduce influence functions for efficient data attribution in learning-based
23 control, enabling practitioners to identify influential training trajectories in seconds rather than hours.
24 Specifically:

- 25 • **IF1 (Dynamics-level):** Estimates how removing a trajectory affects model predictive accu-
26 racy—critical for ensuring accurate state estimation in safety-critical regions.
- 27 • **IF2 (Control-level):** Quantifies impact on LQR control cost, a proxy for tracking error and
28 stability margins. This directly connects training data to controller safety performance.

29 To our knowledge, this is the first work connecting influence functions to control performance,
30 enabling efficient quality assurance for safety-critical learning-based controllers. This extended
31 abstract summarizes key contributions from our recent work Li et al. [2025], with emphasis on safety
32 implications for robotic systems. Full mathematical derivations and extended experimental results
33 appear in Li et al. [2025].

35 2 Problem Setup

36 **System.** Consider a discrete-time system $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}_t$ with state $\mathbf{x}_t \in \mathbb{R}^{n_x}$, control
37 $\mathbf{u}_t \in \mathbb{R}^{n_u}$, and Gaussian noise $\mathbf{w}_t \sim \mathcal{N}(0, \Sigma_w)$.

Learning. We learn a linear model $\hat{\mathbf{x}}_{t+1} = \mathbf{A}_\theta \mathbf{x}_t + \mathbf{B}_\theta \mathbf{u}_t$ from dataset $\mathcal{D} = \{\tau_1, \dots, \tau_N\}$ by minimizing:

$$L(\theta, \mathcal{D}) = \sum_{k=1}^N \mathcal{L}_k(\theta) = \sum_{k=1}^N \sum_{s \in \tau_k} \|\mathbf{x}_{t+1,s} - \Phi_s \theta\|_2^2$$

38 where $\mathcal{L}_k(\theta)$ is the loss on trajectory τ_k .

39 **Control.** Using learned dynamics $(\mathbf{A}_\theta, \mathbf{B}_\theta)$, we design an LQR controller minimizing $J(\theta) =$
40 $\text{Tr}(\mathbf{P}(\theta))$ where $\mathbf{P}(\theta)$ solves the Discrete Algebraic Riccati Equation (DARE). The cost $J(\theta)$
41 represents expected tracking error and state deviation—key safety metrics.

42 **Goal.** Efficiently estimate how removing trajectory τ_k affects: (1) model accuracy $L_{pred}(\theta)$ on test
43 data, and (2) control performance $J(\theta)$, without exhaustive retraining.

44 3 Method Overview

45 3.1 IF1: Influence on Model Accuracy

Removing τ_k changes optimal parameters from $\hat{\theta}$ to $\hat{\theta}_{\setminus k}$. Using influence functions Koh and Liang [2017], we approximate:

$$\Delta \hat{\theta}_k \approx \mathbf{H}_{\hat{\theta}}^{-1} \nabla_{\theta} \mathcal{L}_k(\hat{\theta})$$

46 where $\mathbf{H}_{\hat{\theta}} = \nabla_{\theta}^2 L(\hat{\theta}, \mathcal{D})$ is the Hessian. The influence on test set loss is:

$$IF1(\tau_k, L_{pred}) = (\nabla_{\theta} \mathcal{L}_k(\hat{\theta}))^T \mathbf{H}_{\hat{\theta}}^{-1} \nabla_{\theta} L_{pred}(\hat{\theta}) \quad (1)$$

47 **Safety relevance:** IF1 identifies trajectories whose removal degrades model accuracy in safety-critical
48 state regions, potentially leading to poor state estimation and unsafe control decisions.

49 3.2 IF2: Influence on Control Performance

50 More critically for safety, we quantify how τ_k affects LQR cost:

$$IF2(\tau_k, J) = (\nabla_{\theta} \mathcal{L}_k(\hat{\theta}))^T \mathbf{H}_{\hat{\theta}}^{-1} \nabla_{\theta} J(\hat{\theta}) \quad (2)$$

The key challenge is computing $\nabla_{\theta} J(\hat{\theta})$, which requires tracing sensitivities through the DARE solution. For each parameter θ_m , we compute $S_m = \frac{\partial \mathbf{P}}{\partial \theta_m}$ by solving a Lyapunov equation:

$$S_m - \mathbf{A}_{CL}^T S_m \mathbf{A}_{CL} = -\frac{\partial \mathcal{R}}{\partial \theta_m}$$

51 where $\mathcal{R}(\mathbf{P}, \theta) = 0$ is the DARE and $\mathbf{A}_{CL} = \mathbf{A}_{\hat{\theta}} - \mathbf{B}_{\hat{\theta}} \mathbf{K}(\hat{\theta})$ is the closed-loop system matrix.

52 This yields $\nabla_{\theta} J(\hat{\theta}) = [\text{Tr}(S_1), \dots, \text{Tr}(S_p)]^T$.

53 **Safety relevance:** IF2 identifies trajectories that degrade control performance—manifesting as
54 increased tracking errors, larger state deviations, or reduced stability margins in deployed systems.

55 4 Empirical Validation

56 **Systems.** We validate on two linear systems: Single-link manipulator analogue (2 states, 1 input), Two-
57 link manipulator analogue (4 states, 2 inputs).

58 **Setup.** Training sets: 30-50 trajectories of 25-30 steps. Test sets: 20 trajectories. LQR weights:
59 $\mathbf{Q}_c = \mathbf{I}$, $\mathbf{R}_c = 0.1\mathbf{I}$. For each trajectory τ_k , we computed IF1/IF2 predictions and compared against
60 ground truth via explicit retraining on $\mathcal{D} \setminus \{\tau_k\}$.

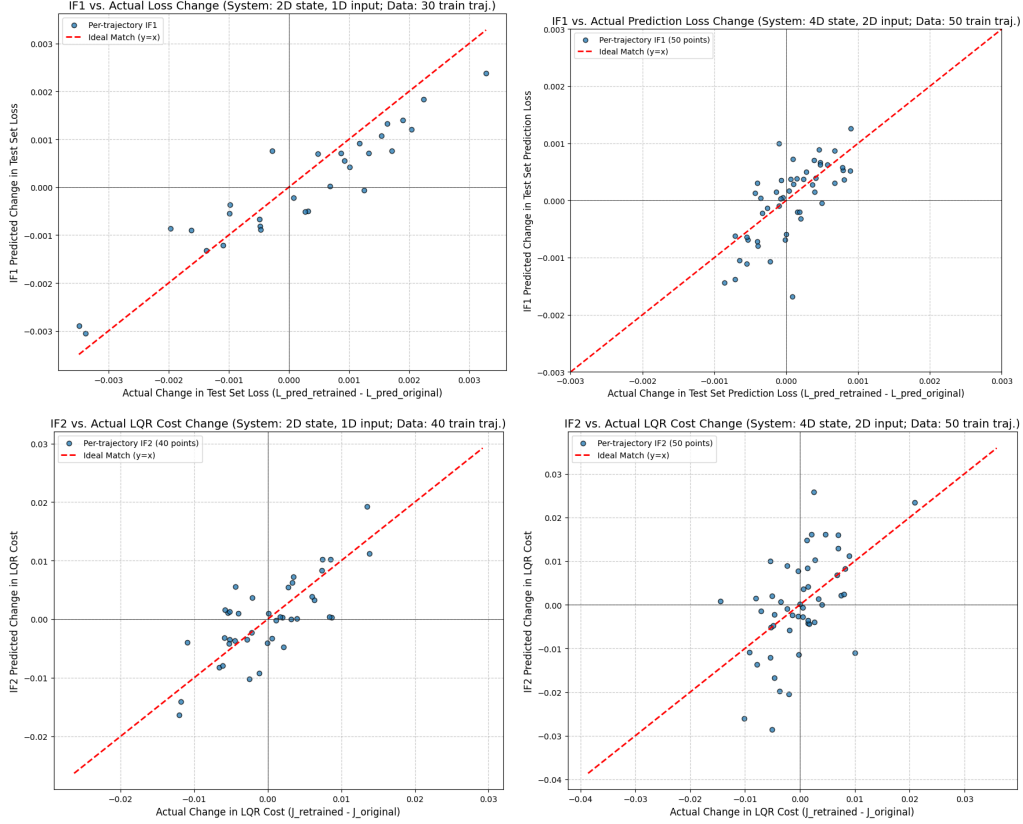


Figure 1: **Validation of influence functions.** Top: IF1 predictions vs. actual changes in test set prediction loss for S1 (left, $r = 0.93$) and S2 (right, $r = 0.69$). Bottom: IF2 predictions vs. actual changes in LQR control cost for S1 (left, $r = 0.71$) and S2 (right, $r = 0.64$). Strong positive correlations demonstrate that influence functions efficiently identify training trajectories that impact controller reliability.

Results. Figure 1 demonstrates strong agreement between influence predictions and ground truth. Critically, IF2 successfully identifies trajectories that affect downstream control performance—enabling efficient quality assurance for safety-critical applications. The moderate scatter in higher-dimensional S2 reflects approximation error from first-order Taylor expansion; second-order corrections could improve accuracy (future work).

Computational Efficiency. Computing IF1/IF2 for all $N = 50$ trajectories required <1 second, versus 8.3 minutes for exhaustive LOO retraining—a $500\times$ speedup.

5 Implications for Safe Robotic Control

5.1 Quality Assurance Applications

Our framework enables several safety-critical capabilities:

Corrupted Data Detection. Trajectories with anomalously high IF2 scores may indicate corrupted or adversarial data that degrades controller performance through systematic model bias. Such data could arise from sensor failures, communication errors, or malicious attacks. Practitioners can flag these for manual review and potential removal to improve controller safety margins.

Critical Data Identification. Trajectories with high positive IF2 influence are critical for maintaining performance in key operating regions. Their removal would increase tracking error and potentially compromise stability—these should be preserved, verified for accuracy, and prioritized during data storage and backup procedures.

79 **Coverage Verification.** Low-influence regions may indicate under-represented operating regimes
80 where the learned model lacks sufficient training support. In safety-critical applications (autonomous
81 driving near obstacles, surgical robots near anatomical boundaries), ensuring adequate data coverage
82 across all operational conditions is essential for reliable deployment.

83 **Active Data Collection.** IF scores can guide additional data gathering in regions where current
84 training data provides weak support for safe control. This enables targeted data acquisition strategies
85 that efficiently improve controller robustness in underrepresented or high-risk scenarios.

86 5.2 Limitations and Future Work

87 **Linear Dynamics and LQR Control.** The current framework is restricted to linear systems and LQR
88 control. Extension to nonlinear dynamics (neural network models) and nonlinear controllers (MPC,
89 iLQR) requires tractable gradient computation through nonlinear Riccati equations and handling
90 influence propagation through neural network parameters.

91 **Stochastic Modeling.** Our formulation treats process noise w_t as given but does not model how
92 training data influences noise covariance estimation. Incorporating influence functions for jointly
93 learned dynamics and noise models would enable more complete data attribution in stochastic control
94 settings.

95 **End-to-End Data-Driven Control.** Modern methods like Data-Enabled Predictive Control (DeePC)
96 and Direct Policy Control (DPC) bypass explicit system identification. Extending influence functions
97 to these frameworks would require tracing data influence directly through Hankel matrices or policy
98 parameterizations to closed-loop performance.

99 **Hessian Computation.** Computing the full Hessian $H_{\hat{\theta}}$ becomes prohibitive for high-dimensional
100 systems. Future work should explore Gauss-Newton or Fisher information approximations, implicit
101 Hessian-vector products via automatic differentiation, and Hessian-free methods using conjugate
102 gradient or Neumann series to enable scaling to large neural network models.

103 **Approximation Error.** First-order Taylor approximation may be inaccurate for large parameter
104 perturbations. Deriving theoretical error bounds or developing adaptive second-order corrections
105 would strengthen reliability guarantees.

106 6 Conclusion

107 We introduced influence functions (IF1, IF2) for efficient data attribution in learning-based control
108 systems. IF1 quantifies how training trajectories affect model predictive accuracy, while IF2 represents
109 the first method to efficiently trace training data influence through system identification to control
110 performance—directly connecting individual trajectories to LQR cost and controller reliability. This
111 framework enables quality assurance without exhaustive retraining. It provides practitioners with a
112 computationally tractable tool for identifying influential data in safety-critical applications.

113 References

- 114 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
115 Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 116 Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the sample complexity
117 of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20:633–679, 2020.
- 118 Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long
119 tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891,
120 2020.
- 121 Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit
122 Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. Studying large language model generalization
123 with influence functions. *arXiv preprint arXiv:2308.03296*, 2023.
- 124 Zayd Hammoudeh and Daniel Lowd. Training data influence analysis and estimation: A survey.
125 In *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in*
126 *Databases*, pages 157–172. Springer, 2022.

- 127 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
128 *International Conference on Machine Learning*, pages 1885–1894. PMLR, 2017.
- 129 Jesse Levinson, Jake Askeland, Jan Becker, Jennifer Dolson, David Held, Soeren Kammel, J Zico
130 Kolter, Dirk Langer, Oliver Pink, Vaughan Pratt, et al. Towards fully autonomous driving: Systems
131 and algorithms. *IEEE Intelligent Vehicles Symposium*, pages 163–168, 2011.
- 132 Jiachen Li, Shihao Li, Jiamin Xu, Soovadeep Bakshi, and Dongmei Chen. Influence functions for
133 data attribution in linear system identification and lqr control, 2025. URL <https://arxiv.org/abs/2506.11293>.
134
- 135 Azad Shademan, Ryan S Decker, Justin D Opfermann, Simon Leonard, Axel Krieger, and Peter CW
136 Kim. Supervised autonomous robotic soft tissue surgery. *Science Translational Medicine*, 8(337):
137 337ra64–337ra64, 2016.