

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser^{*,†}, Alexandru Tifrea^{*,†}, Michael Aerni[†], Reinhard Heckel^{°,§}, Fanny Yang[†]
 ETH Zurich[†], Rice University[°], TU Munich[§]



PHENOMENON 1: DOUBLE DESCENT

Observed empirically for neural networks and theoretically e.g. for highly overparameterized ($d \gg n$) linear models [1].

- ▶ Regularization does not improve generalization, compared to interpolating the training data.
- ▶ Overparameterization implicitly controls the variance. → Regularization is **redundant**.

PHENOMENON 2: ROBUST OVERFITTING

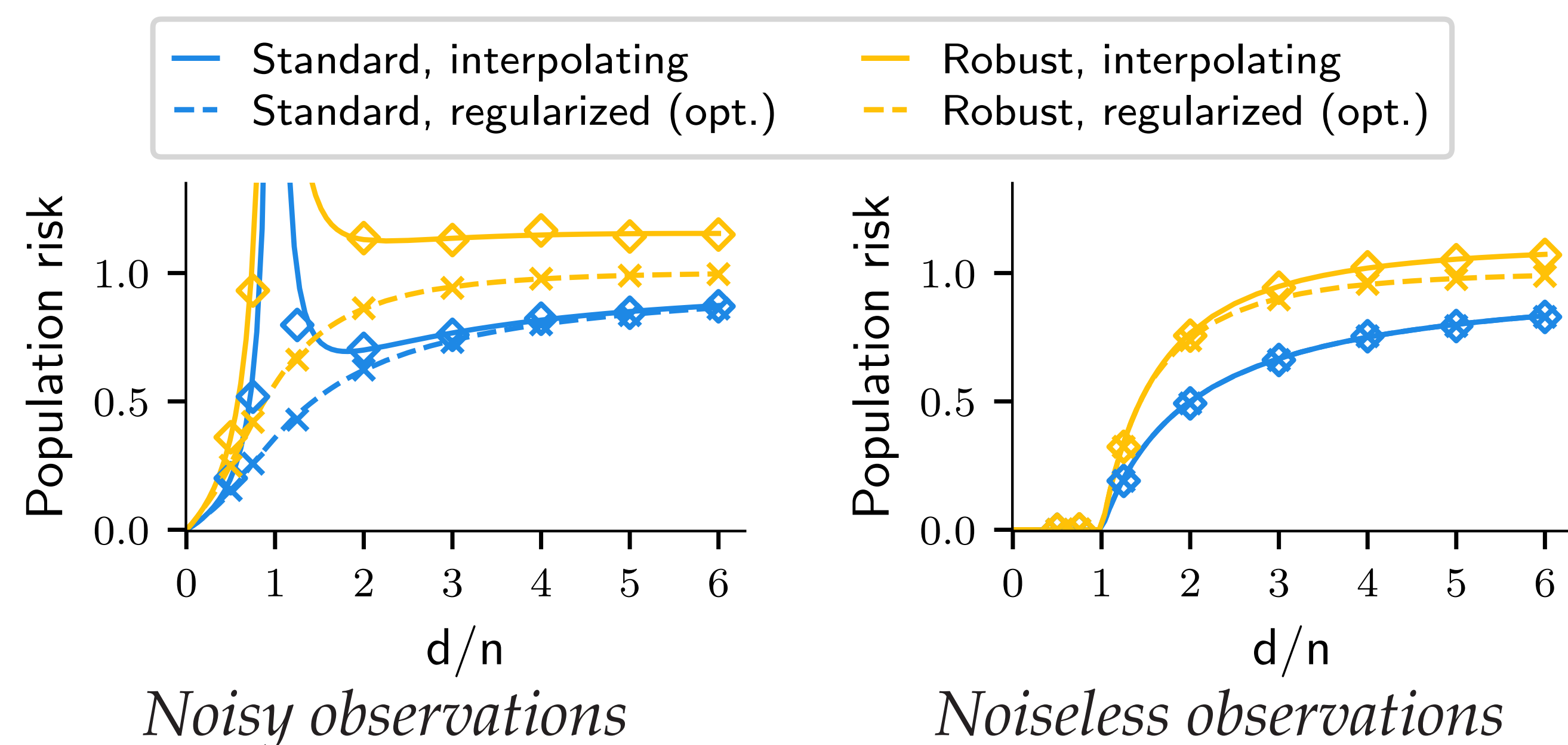
Observed empirically for neural networks on image data [2].

- ▶ *Robust* generalization benefits greatly from regularization.
- ▶ Prior work has attributed this phenomenon to:
 - noise in the training data
 - non-smooth predictors

Does robust overfitting occur on noiseless data?
 Does this provably happen even for linear models?

ROBUST LINEAR REGRESSION

Ridge regularization **avoids the min-norm interpolator**.



- ▶ The lowest robust risks are not obtained by the min-norm interpolators, but by the regularized estimators.
 - holds true even for **noiseless data!**

ROBUST LINEAR CLASSIFICATION

▶ Evaluation with the **robust risk** wrt ℓ_∞ -perturbations:

$$\mathbf{R}_\epsilon(\theta) := \mathbb{E}_{X \sim \mathcal{P}} \max_{\delta \in \mathcal{U}_c(\epsilon)} \mathbb{1}_{\text{sgn}(\langle \theta, X + \delta \rangle) \neq \text{sgn}(\langle \theta^*, X \rangle)}$$

▶ We use adversarial training to obtain a robust estimator:

$$\hat{\theta}_\lambda := \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\delta \in \mathcal{U}_c(\epsilon)} \ell(\langle \theta, x_i + \delta \rangle y_i) + \lambda \|\theta\|_2^2$$

▶ For $\lambda \rightarrow 0$, it maximizes the robust margin of the data:

$$\hat{\theta}_0 := \arg \min_{\theta} \|\theta\|_2 \text{ such that for all } i, \max_{\delta \in \mathcal{U}_c(\epsilon)} y_i \langle \theta, x_i + \delta \rangle \geq 1.$$

THEORETICAL RESULT FOR CLASSIFICATION

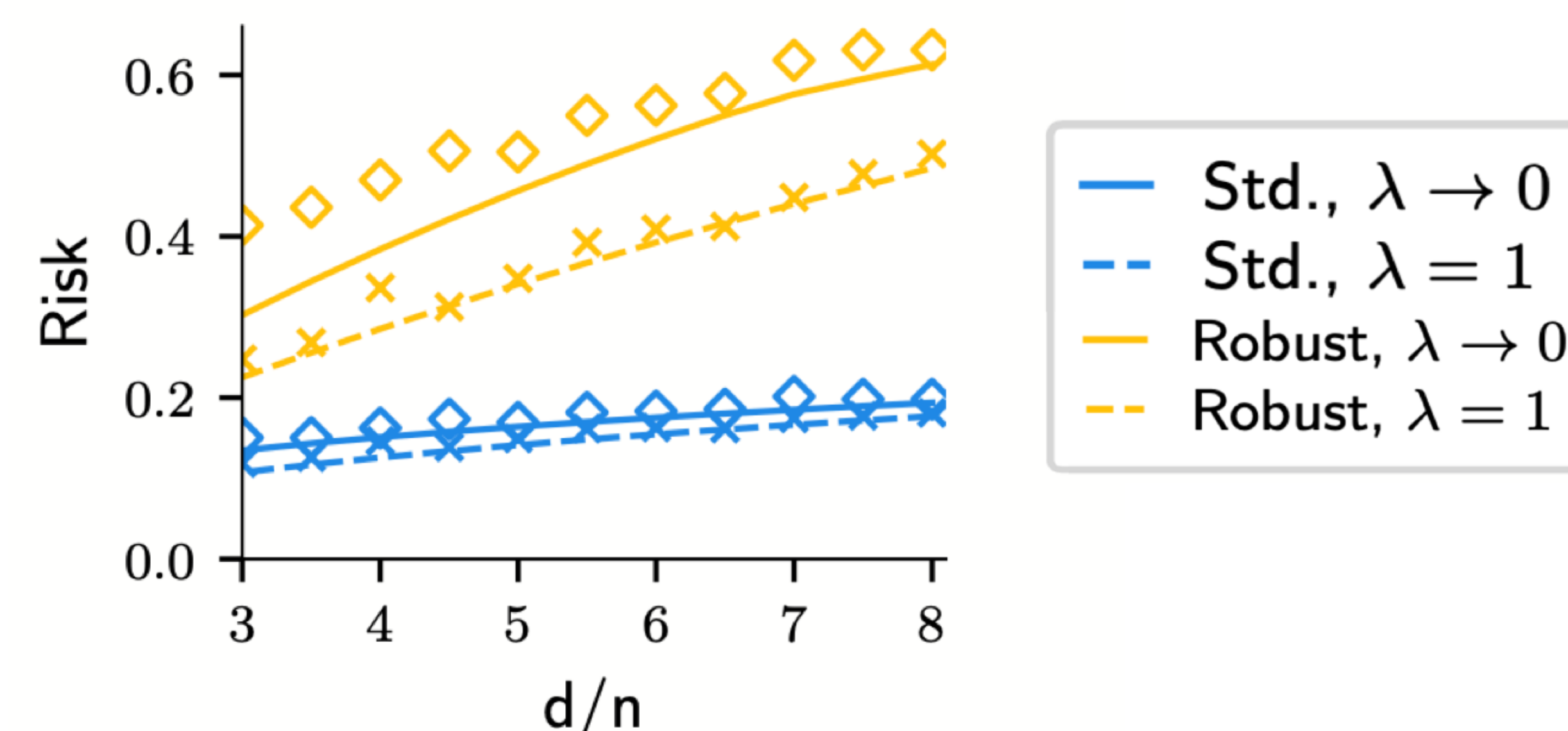
Problem setting:

- ▶ Data model: covariates $x \sim \mathcal{N}(0, I_d)$, deterministic labels given by $y = \text{sgn}(\langle \theta^*, x \rangle) \in \{-1, +1\}$. → **Noiseless data!**
- ▶ We consider **linear classifiers** trained with the logistic loss.

Theorem. For a sparse ground truth, we derive the limit $\mathcal{R}_\lambda(\epsilon, \gamma)$ of the robust risk as $d, n \rightarrow \infty$ and $d/n \rightarrow \gamma$:

$$\mathbf{R}_\epsilon(\hat{\theta}_\lambda) \xrightarrow{\text{prob}} \mathcal{R}_\lambda(\epsilon, \gamma)$$

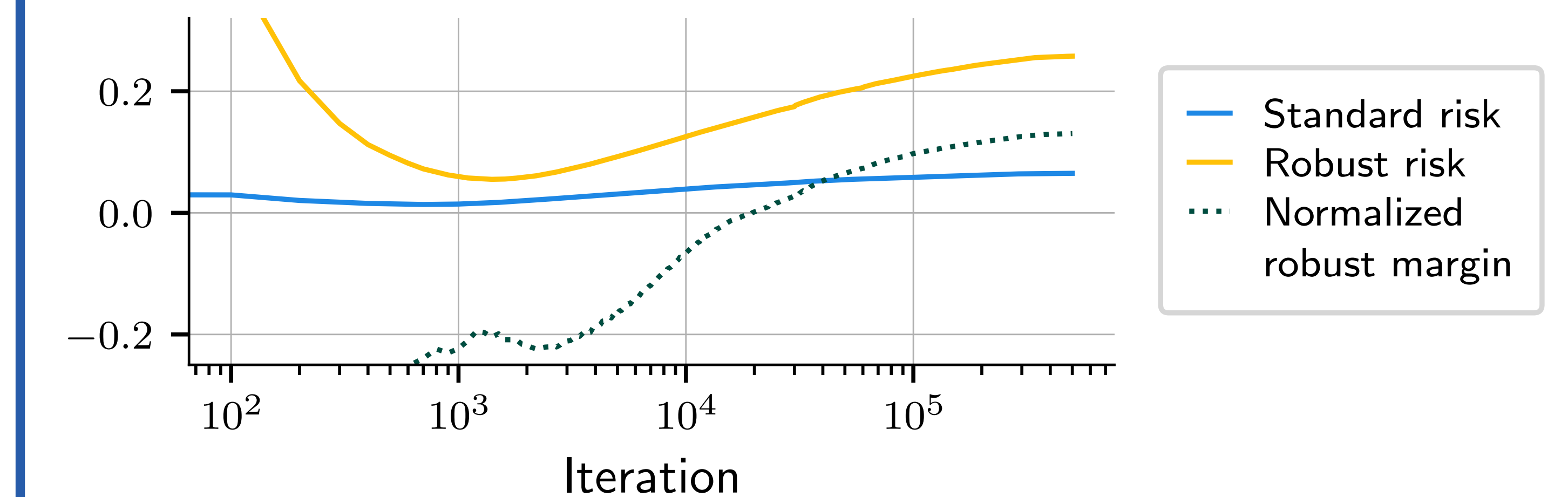
In particular, for some $\lambda_{\text{opt}} > 0$: $\underbrace{\mathcal{R}_{\lambda_{\text{opt}}}(\epsilon, \gamma)}_{\text{regularized}} < \underbrace{\lim_{\lambda \rightarrow 0} \mathcal{R}_\lambda(\epsilon, \gamma)}_{\text{interpolating}}$.



- Lines:** asymptotic risks (theory)
Markers: risks for finite d, n (simulations)
- ▶ Ridge regularization **avoids the max-margin estimator**.

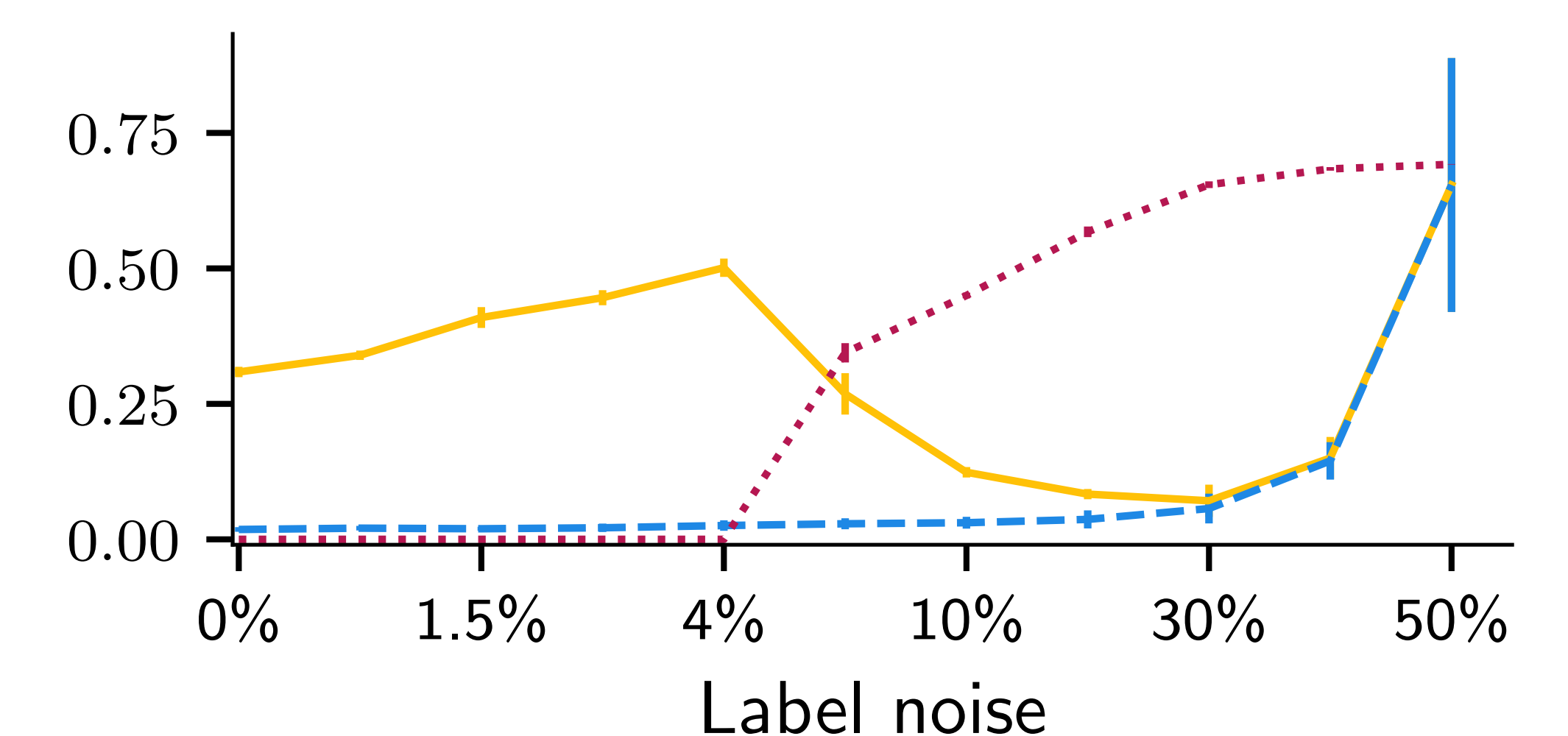
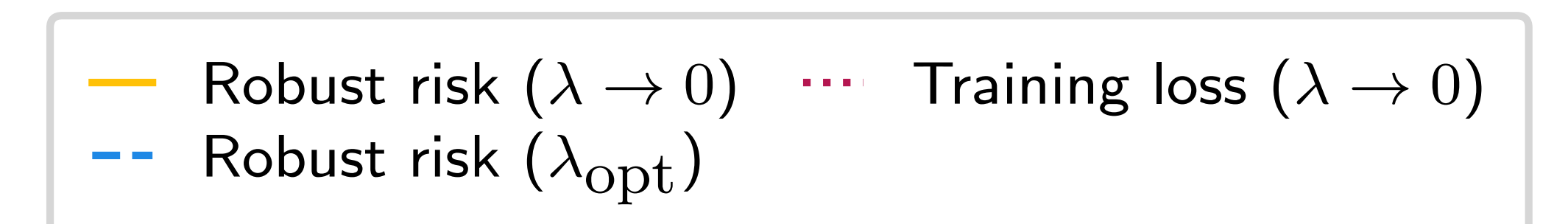
OTHER WAYS TO AVOID $\hat{\theta}_0$

1. Early stopping **avoids the max-margin estimator** and achieves a lower robust risk.



2. Adding artificial label noise prevents a vanishing training loss and **avoids the max-margin estimator**.

Surprising consequence: Smaller robust risk, compared to the max-margin interpolator of the original clean data.



Remark: Regularization still leads to smaller robust risk, even in the presence of noise.

REFERENCES

[1] T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani, "Surprises in high-dimensional ridgeless least squares interpolation," *arXiv:1903.08560*, 2019.
 [2] L. Rice, E. Wong, and Z. Kolter, "Overfitting in adversarially robust deep learning," in *ICML*, 2020.
 [3] C. Thrampoulidis, S. Oymak, and B. Hassibi, "Regularized linear regression: A precise analysis of the estimation error," in *COLT*, 2015.