# Understanding The Failure Modes of Out-of-Distribution Generalization

Vaishnavh Nagarajan<sup>i</sup> Anders Andreassen<sup>2</sup> Behnam Neyshabur<sup>2</sup>

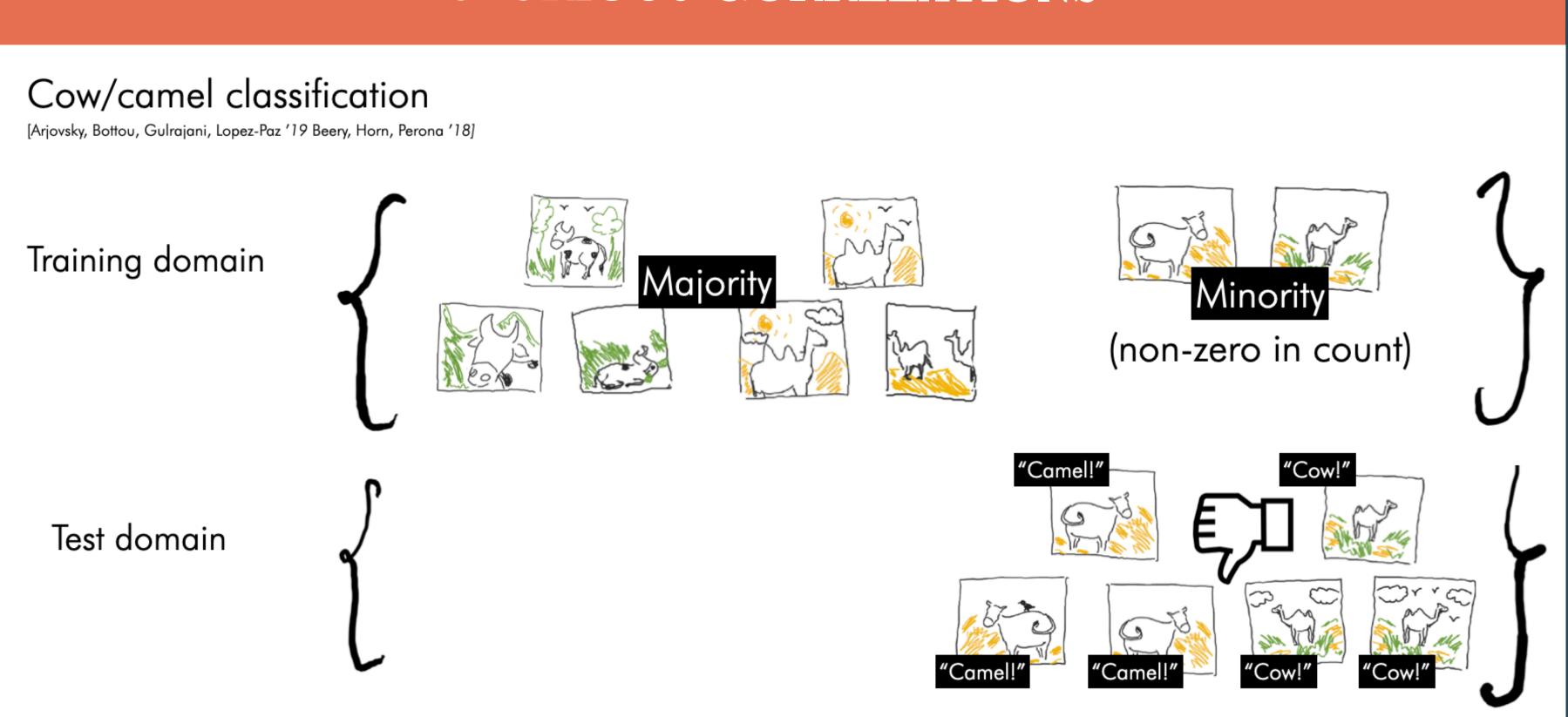
<sup>1</sup>Computer Science Department, Carnegie Mellon University <sup>2</sup>Blueshift, Alphabet

#### HIGH LEVEL MESSAGE

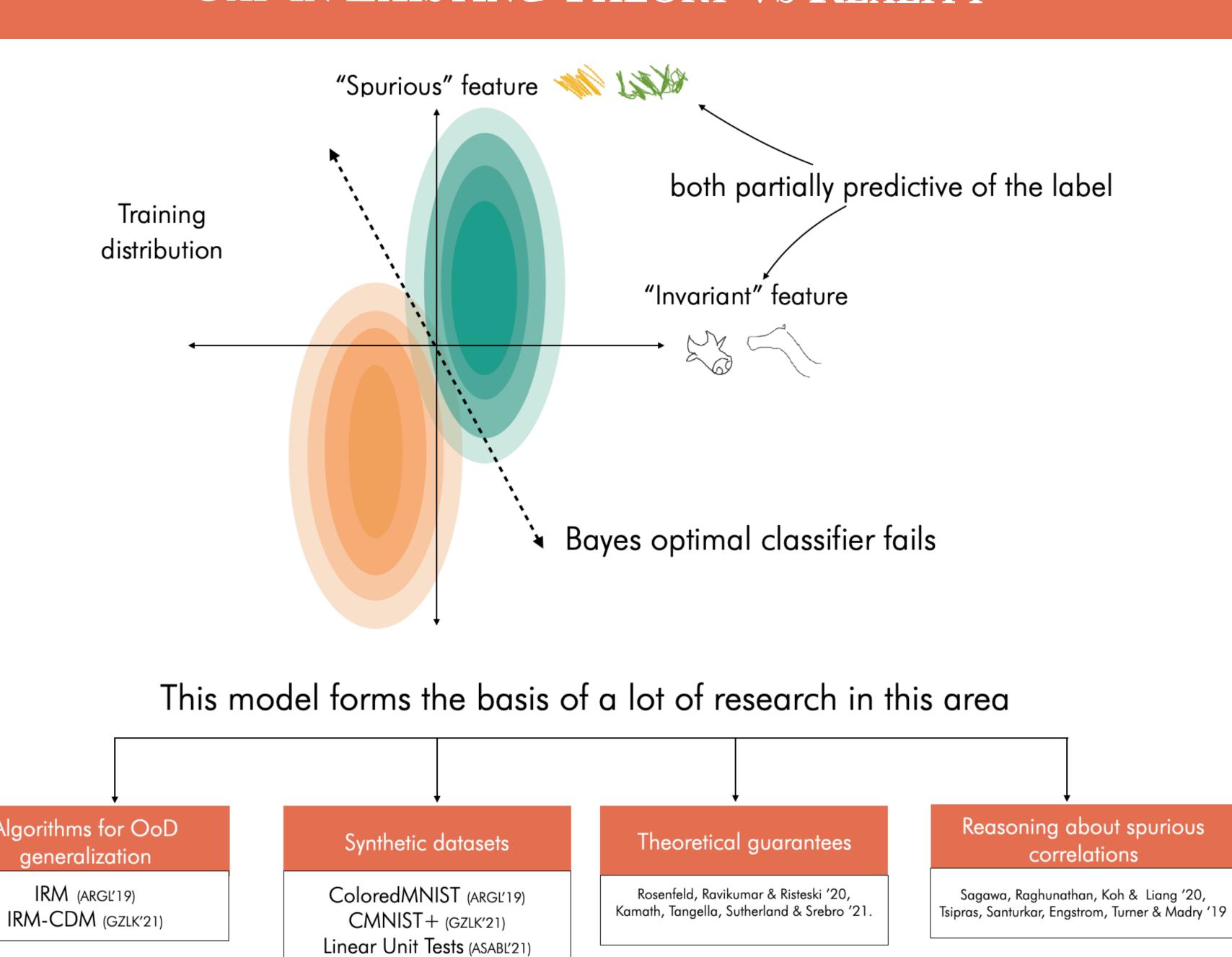
## Why do classifiers rely on spurious correlations?

- 1. Existing theory does not capture the fundamental reasons!
- 2. We theoretically study GD-trained linear classifiers and discover two fundamental failure modes.
- 3. We empirically verify these failure modes in deep learning.

## SPURIOUS CORRELATIONS

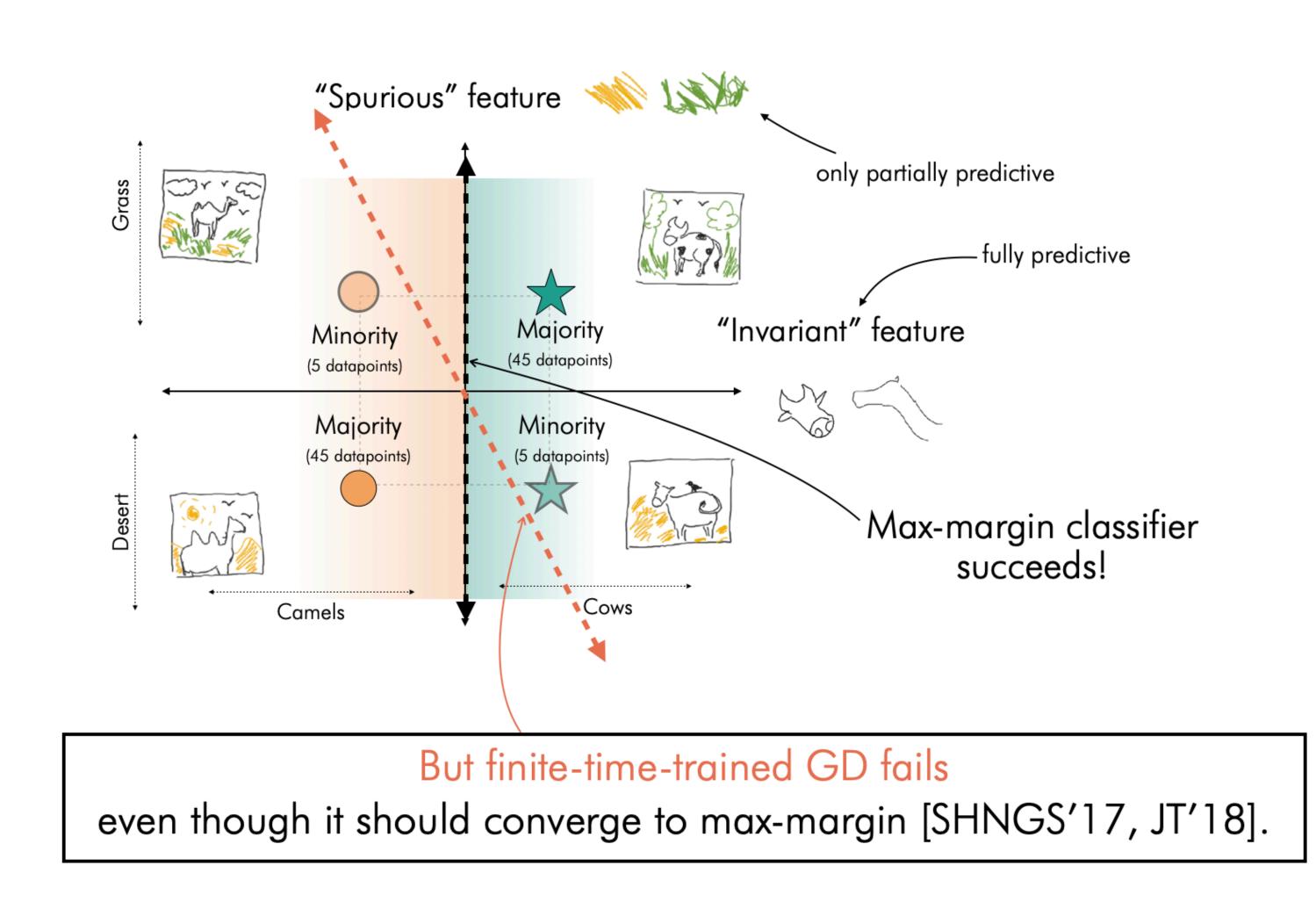


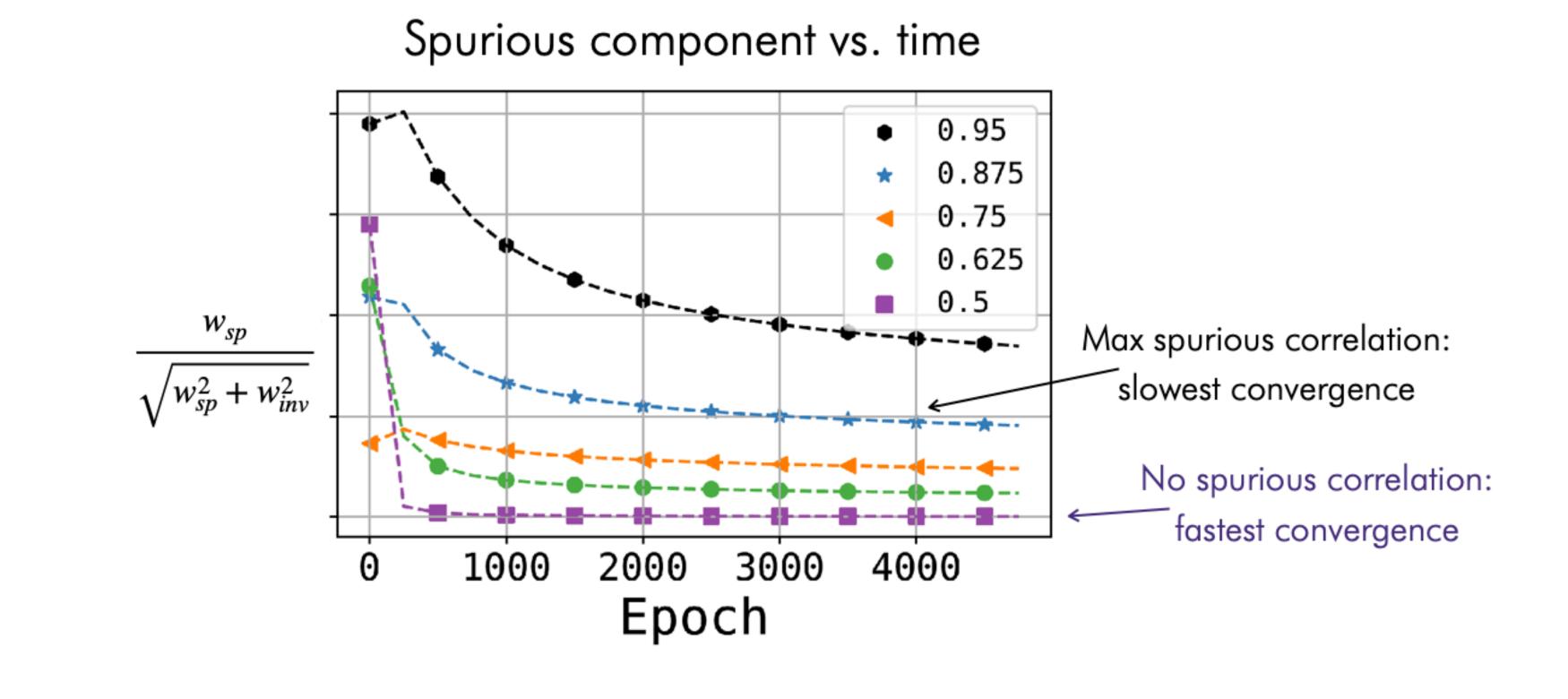
### GAP IN EXISTING THEORY VS REALITY

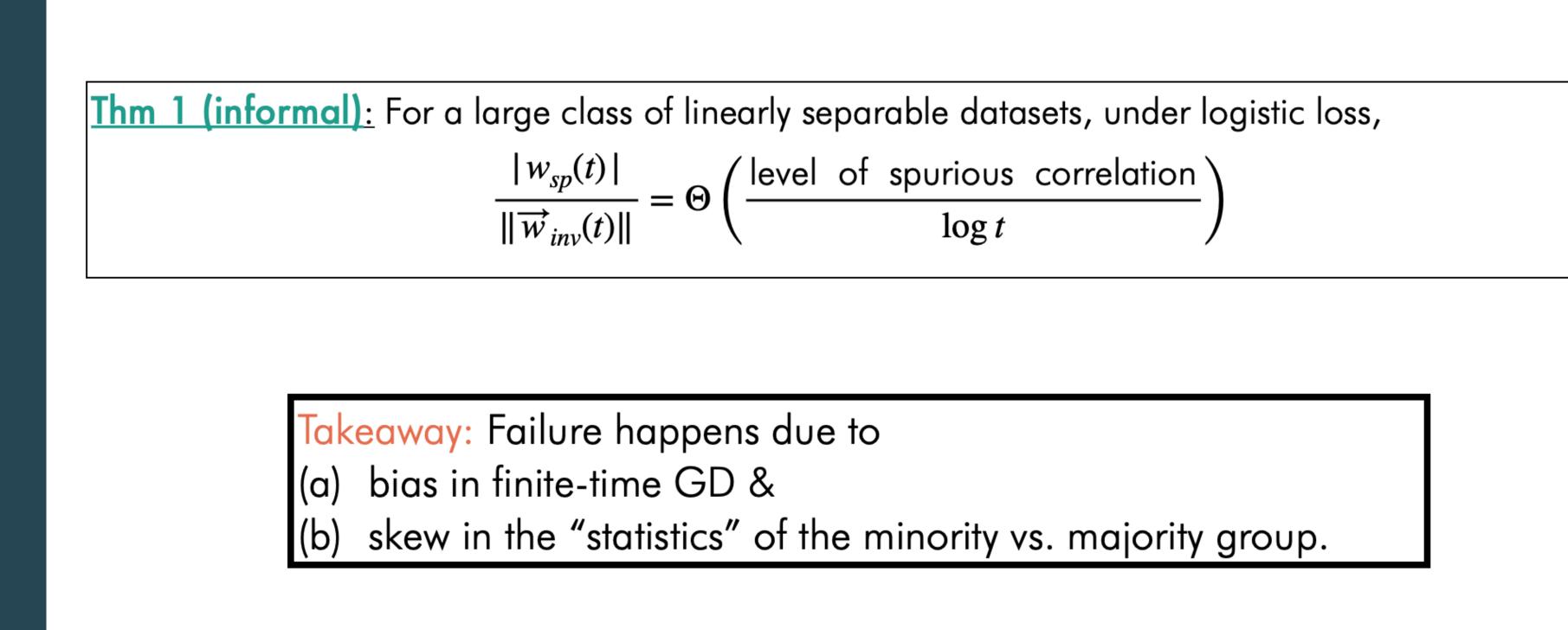


In practice, classifiers use the partially predictive spurious feature, even when the invariant feature is fully predictive!

## FAILURE MODE 1: STATISTICAL SKEWS



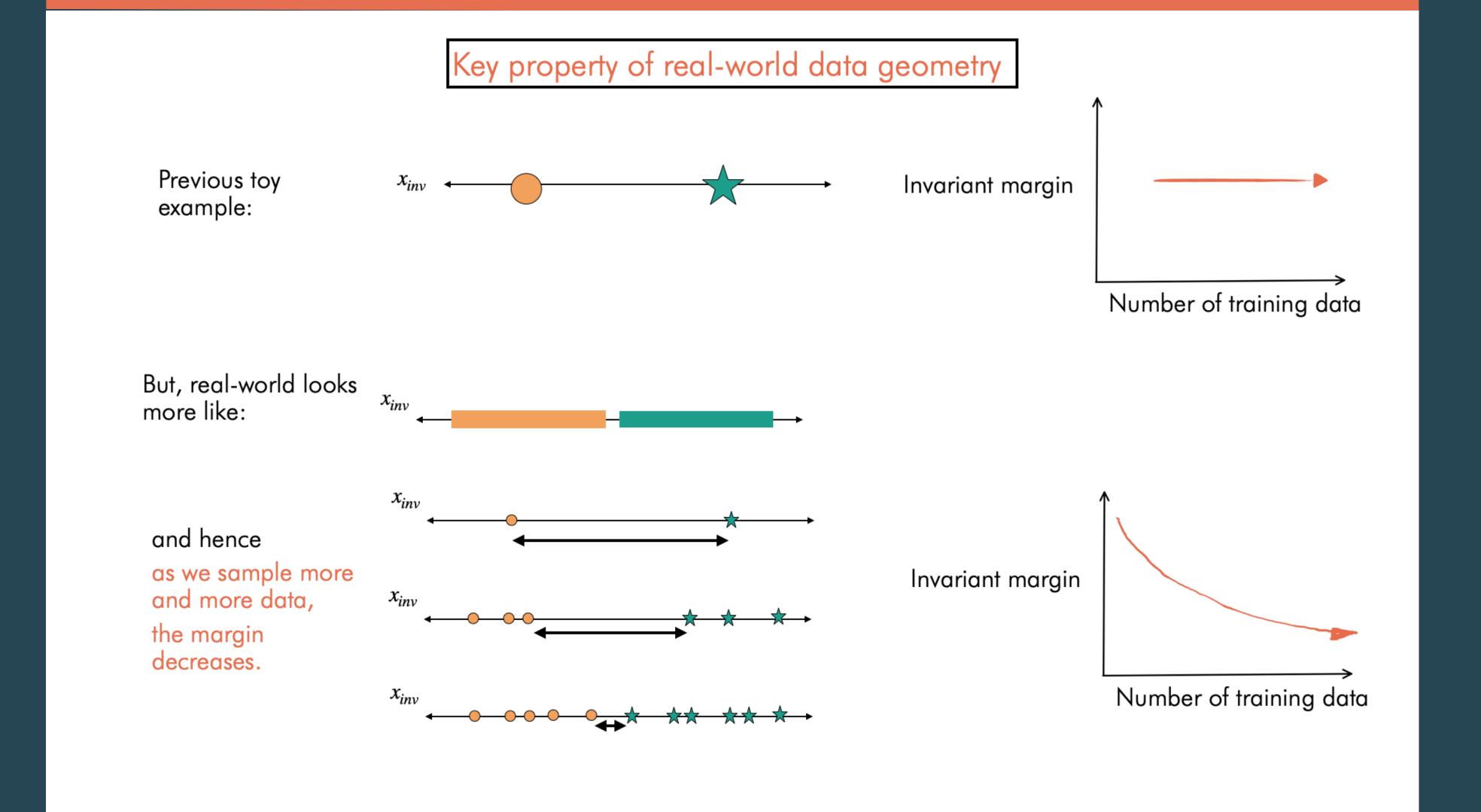


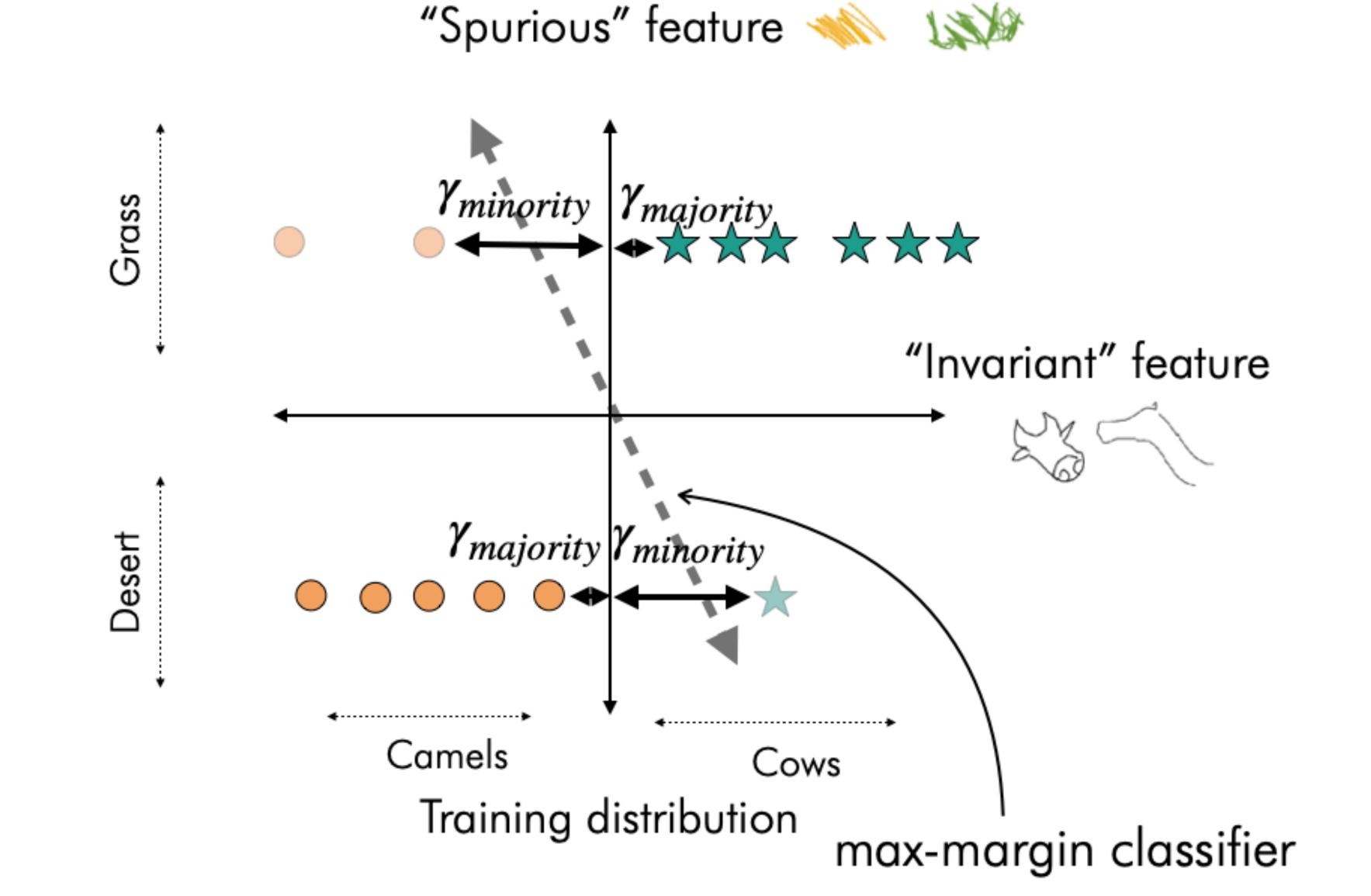


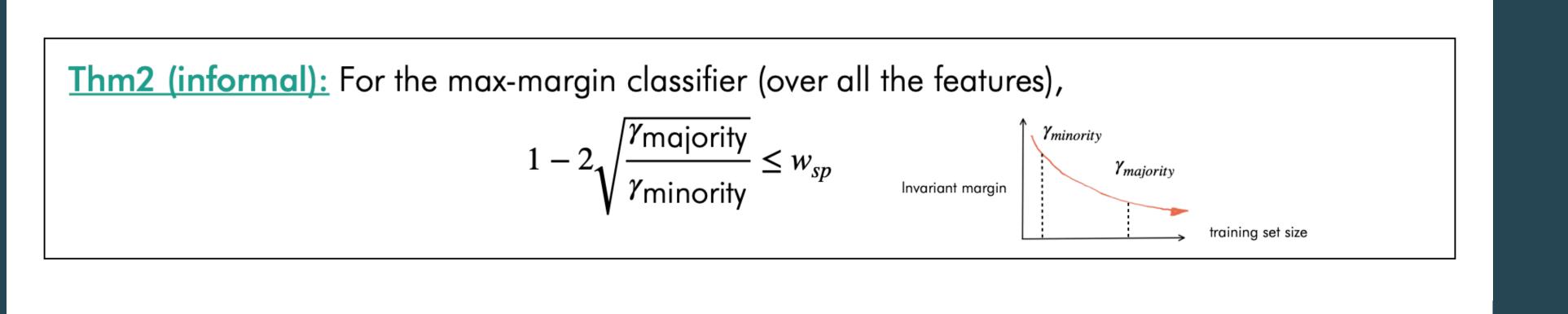
## DOES MAX-MARGIN ALWAYS SUCCEED?

No! We show that when data has non-degenerate geometry, even max-margin classifier can fail...

## FAILURE MODE 2: GEOMETRIC SKEWS



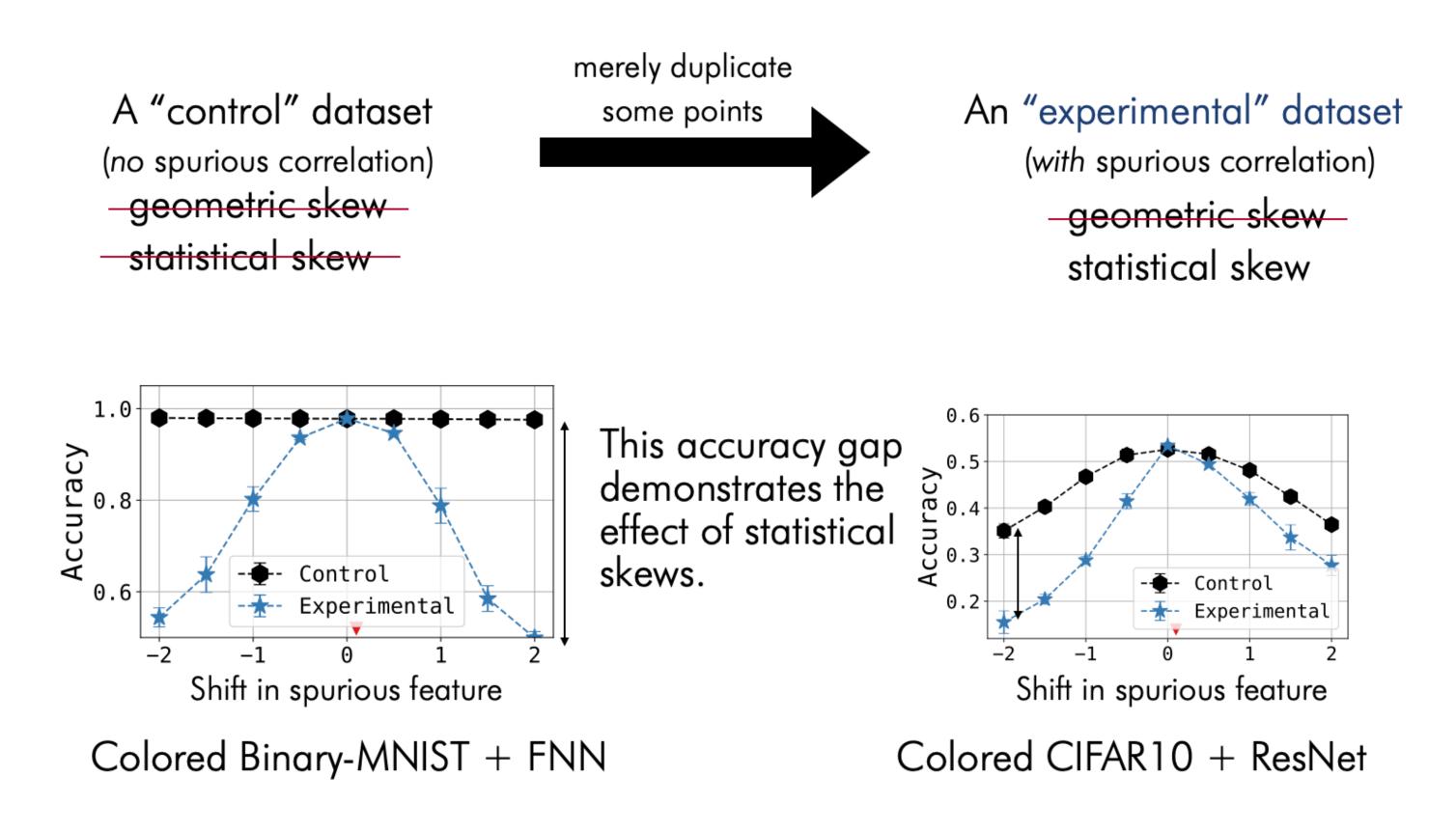




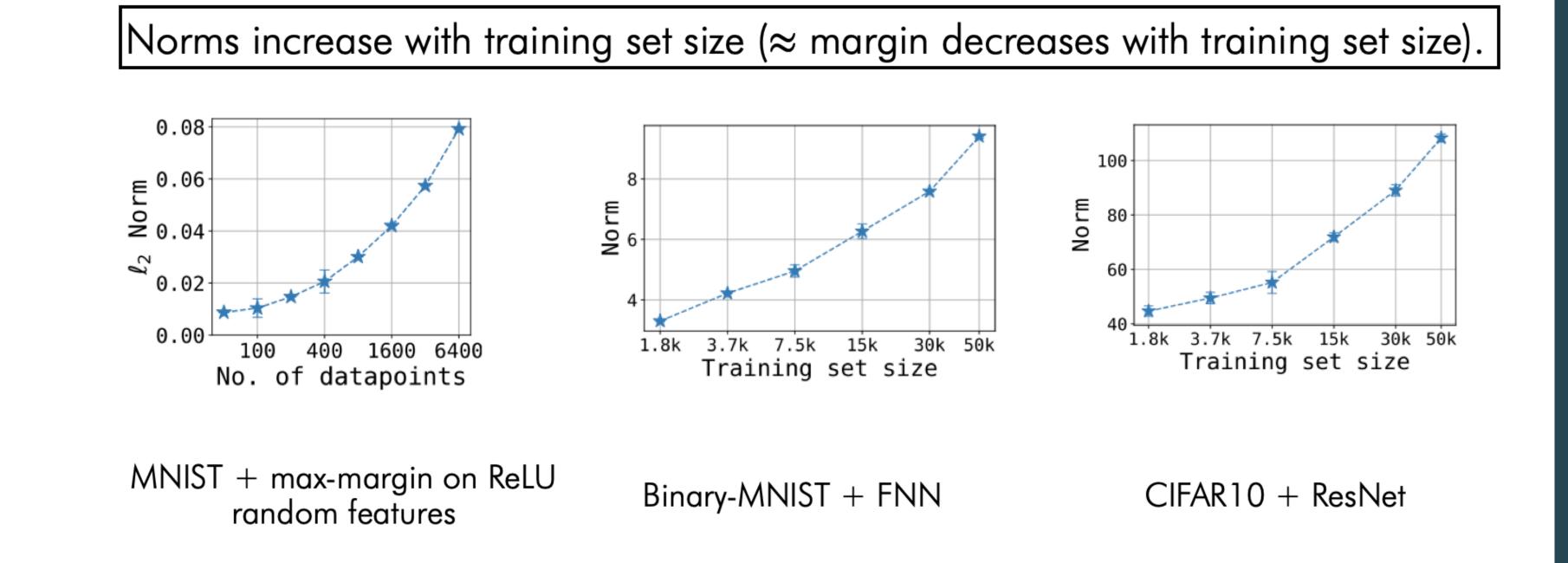
Takeaway: Failure happens because of

- (a) the skew in the geometry of the minority and majority group &
- (b) margin-maximizing bias.

## EMPIRICAL VERIFICATION: STATISTICAL SKEWS



## EMPIRICAL VERIFICATION: GEOMETRIC SKEWS



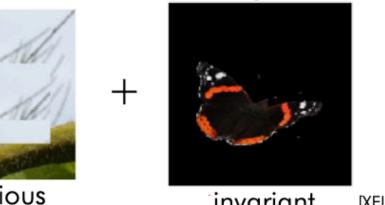
### FUTURE WORK

No one unique way by which classifiers fail!)

- OoD failure modes unique to how deep representations are learned?
- OoD algorithms specialized towards the two skews?
- Characterize real-world "spurious features" vs. "invariant features"?







## REFERENCES

- Arjovsky, Bottou, Gulrajani and Lopez-Paz, "Invariant Risk Minimization", 2019.
- Guo, Zhang, Lio, Kiciman, "Out-of-distribution prediction with IRM: The limitation and an effective fix", 2021
- Aubin, Slowik, Arjovsky, Bottou, Lopez-Paz, "Linear unit-tests for invariance discovery", CDCI ML workshop NeurIPS 2021
- Soudry, Hoffer, Nacson, Gunasekar, and Srebro. "The implicit bias of gradient descent on separable data". J. Mach. Learn. Res., 19, 2018
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression
- Sagawa, Raghunathan, Koh, and Liang. An investigation of why overparameterization exacerbates spurious correlations. 2020