

The Price of Implicit Bias in Robust ML

Nikolaos Tsilivis
NYU

Natalie Frank
University of Washington

Nati Srebro
TTIC

Julia Kempe
NYU & FAIR, Meta

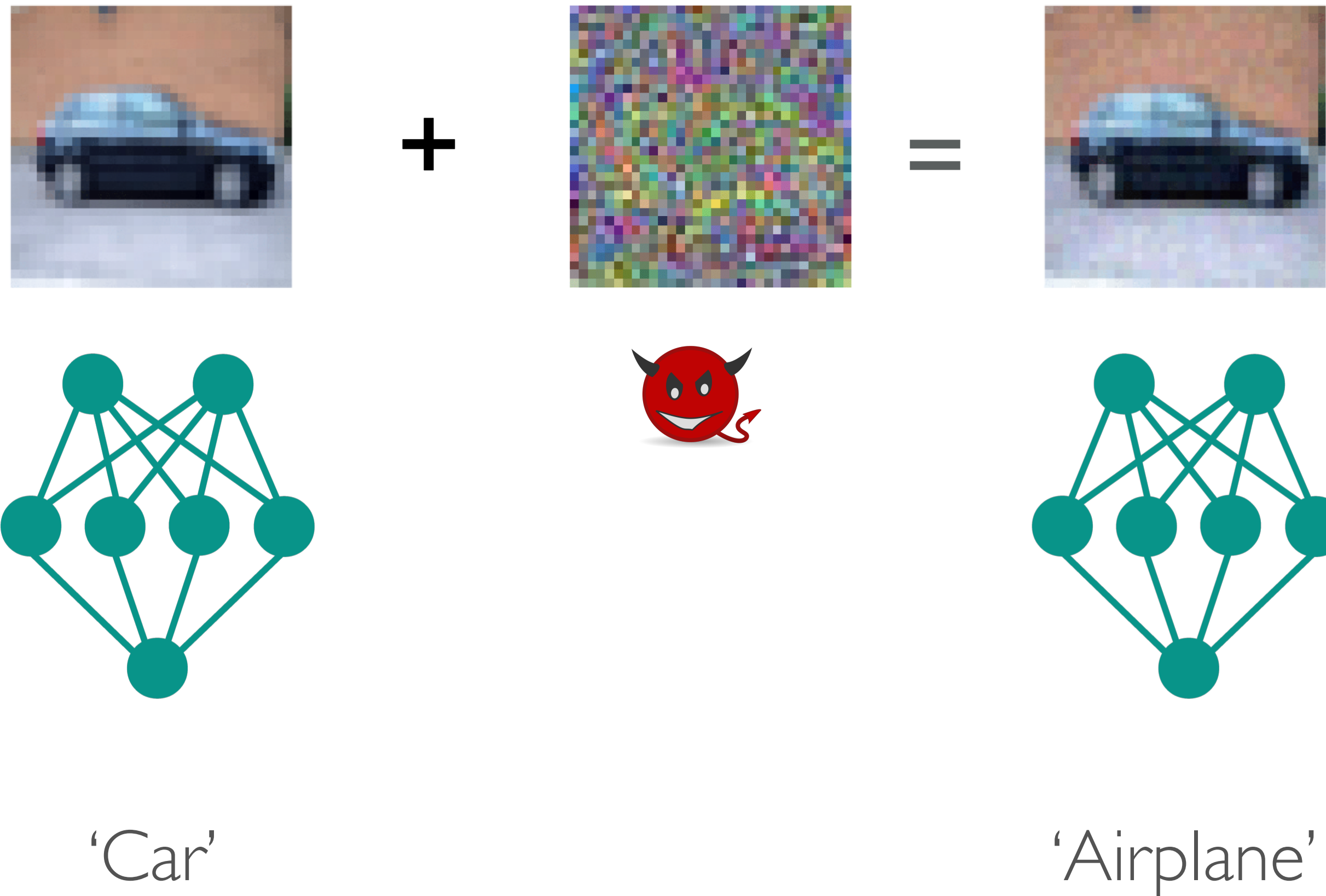
ROBUST SUPERVISED LEARNING

Examples $\{(sample, label)\}$ \longrightarrow Do well in $(sample_new, label_new)$



Test time robustness : Examples $\{(sample, label)\}$ \longrightarrow Do well in *corrupted* $(sample_new', label_new')$

ADVERSARIAL ATTACKS?

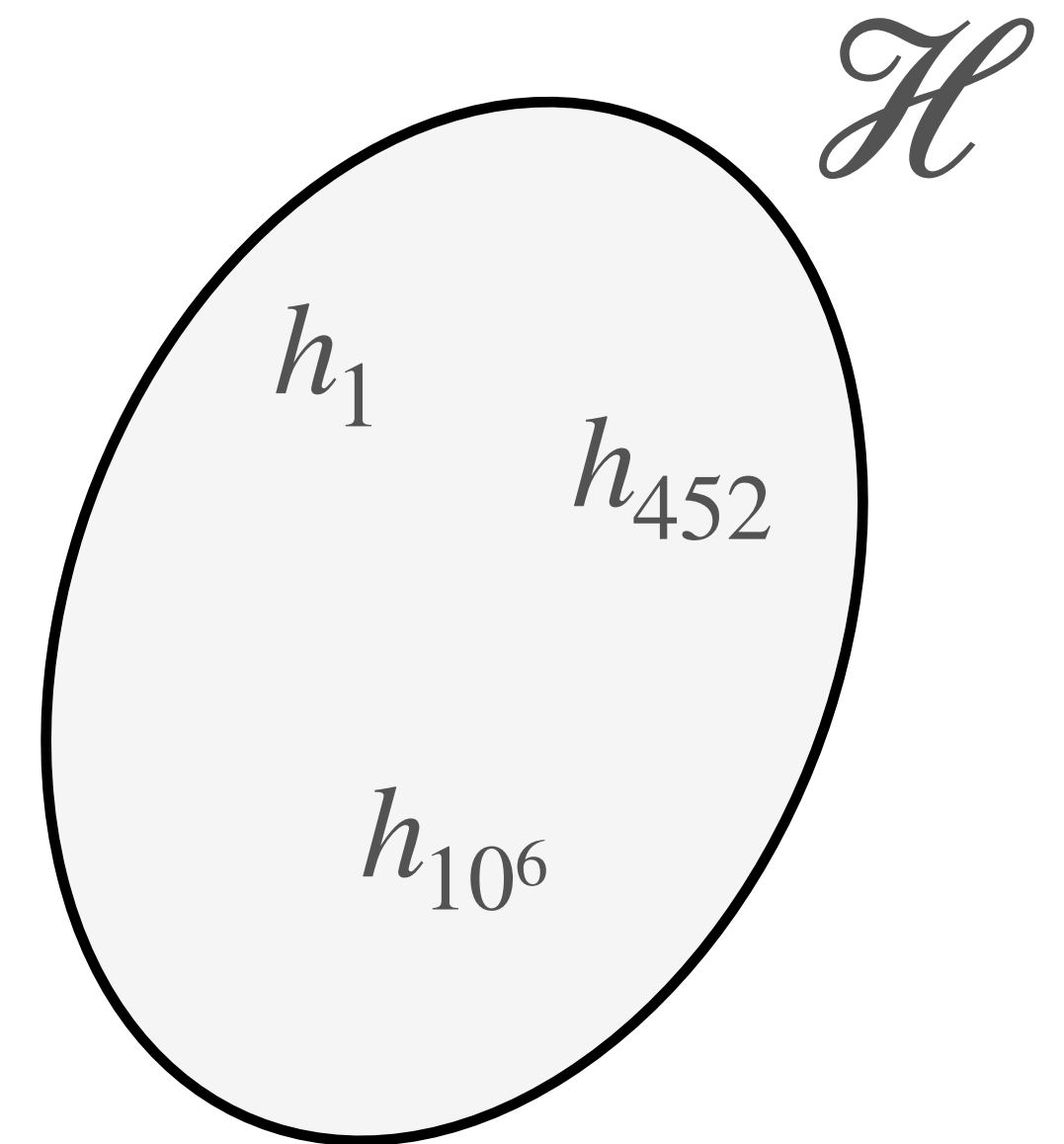


Szegedy et al. Intriguing properties of neural networks. ICLR 2014
Test of Time Award Runner Up, ICLR 2024

ROBUST EMPIRICAL RISK MINIMIZATION

Class of hypotheses \mathcal{H} , dataset $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$, perturbation level ϵ , threat model of ℓ_p balls:

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \mathbf{1} [h(\mathbf{x}'_i) \neq y_i]$$



What is different in Robust ERM???

AGENDA

- Understand what is happening in linear models
 - (Robust) Generalization of linear models
 - ↓↑
 - (Robust) Optimization of linear models
- Suggest what is happening in general

A PRIMER IN LEARNING THEORY (1/2)



Vapnik & Chervonenkis, 1971

Setup: Learn a concept c^* using data $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ from distribution D using class of hypotheses \mathcal{H} . For $h_{ERM} \in \arg \min_{h \in \mathcal{H}} L_S(h)$,

it holds with high probability:

$$L_D(h_{ERM}) \leq \underbrace{\inf_{h \in \mathcal{H}} L_D(h)}_{\text{Approximation error}} + \underbrace{\sqrt{\frac{\text{capacity}(\mathcal{H})}{m}}}_{\text{Estimation error}}$$

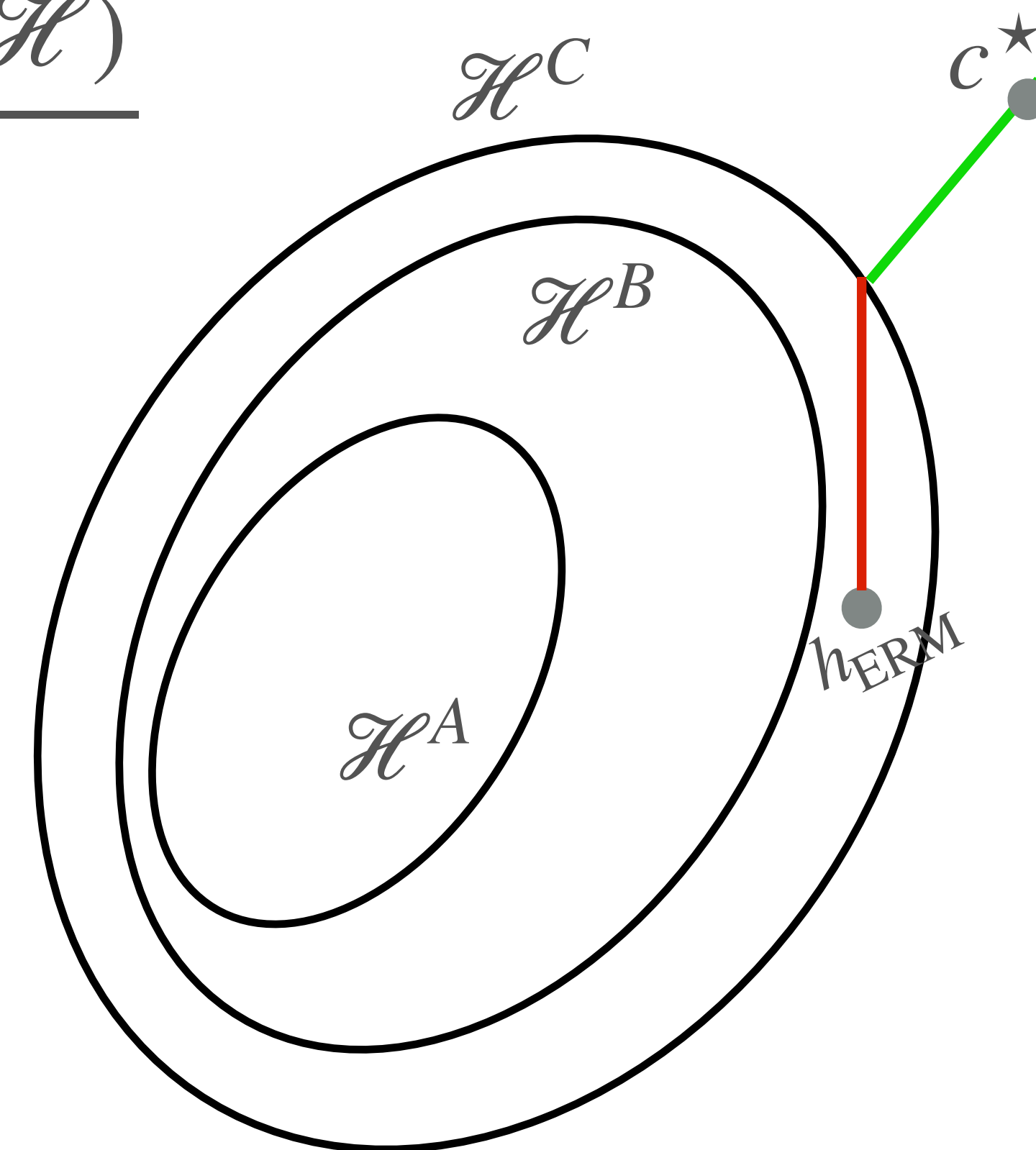
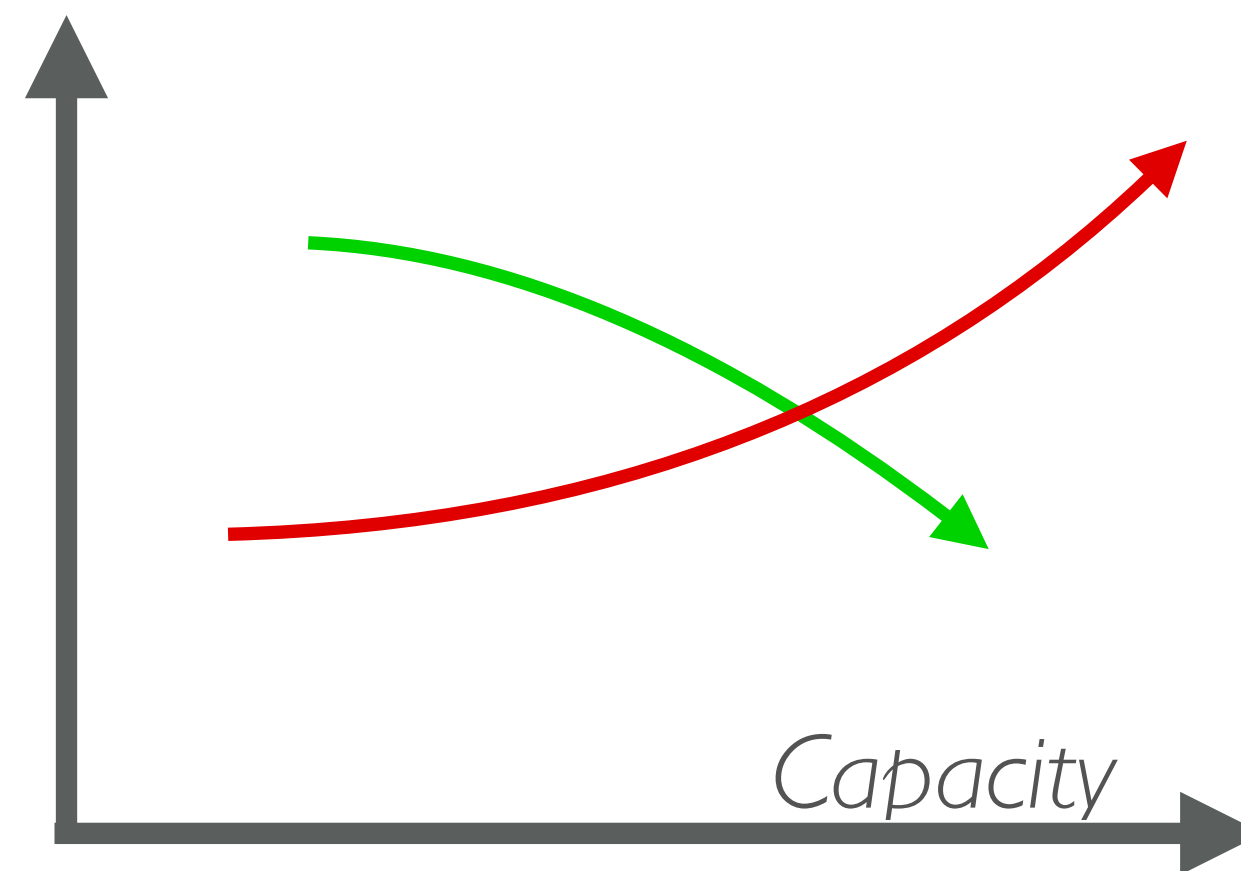
expected error

$$L_D(h) = \mathbb{E}_{(\mathbf{x}, y) \sim D} [l(h(\mathbf{x}), y)]$$

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h(\mathbf{x}_i), y_i)$$

empirical error

Tradeoff between fitness and complexity



A PRIMER IN LEARNING THEORY (2/2)

Setup: Data $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ from distribution D and hypothesis class \mathcal{H} .

Measure of Complexity: VC dimension

Measure of Complexity: **Rademacher complexity**

$$\text{Rad}(\mathcal{H}) = \mathbb{E}_{\sigma_1, \dots, \sigma_m} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i h(\mathbf{x}_i) \right],$$

where $\sigma_i \sim \text{Rad}$ for all $i = 1, \dots, m$.

expected error

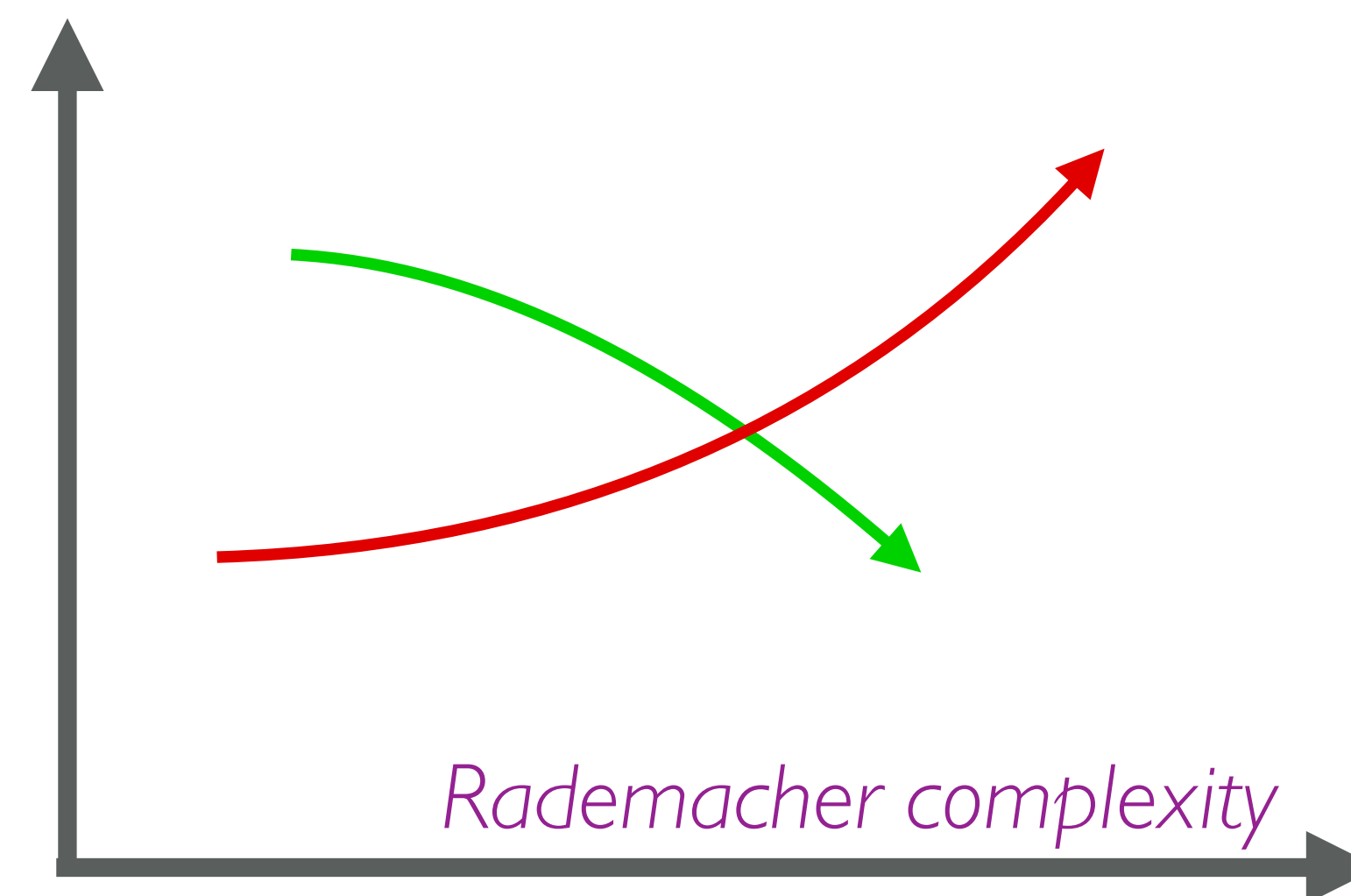
$$L_D(h) = \mathbb{E}_{(\mathbf{x}, y) \sim D} [l(h(\mathbf{x}), y)]$$

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h(\mathbf{x}_i), y_i)$$

empirical error

$$L_D(h_{\text{ERM}}) \leq \inf_{h \in \mathcal{H}} L_D(h) + \text{Rad}(\mathcal{H})$$

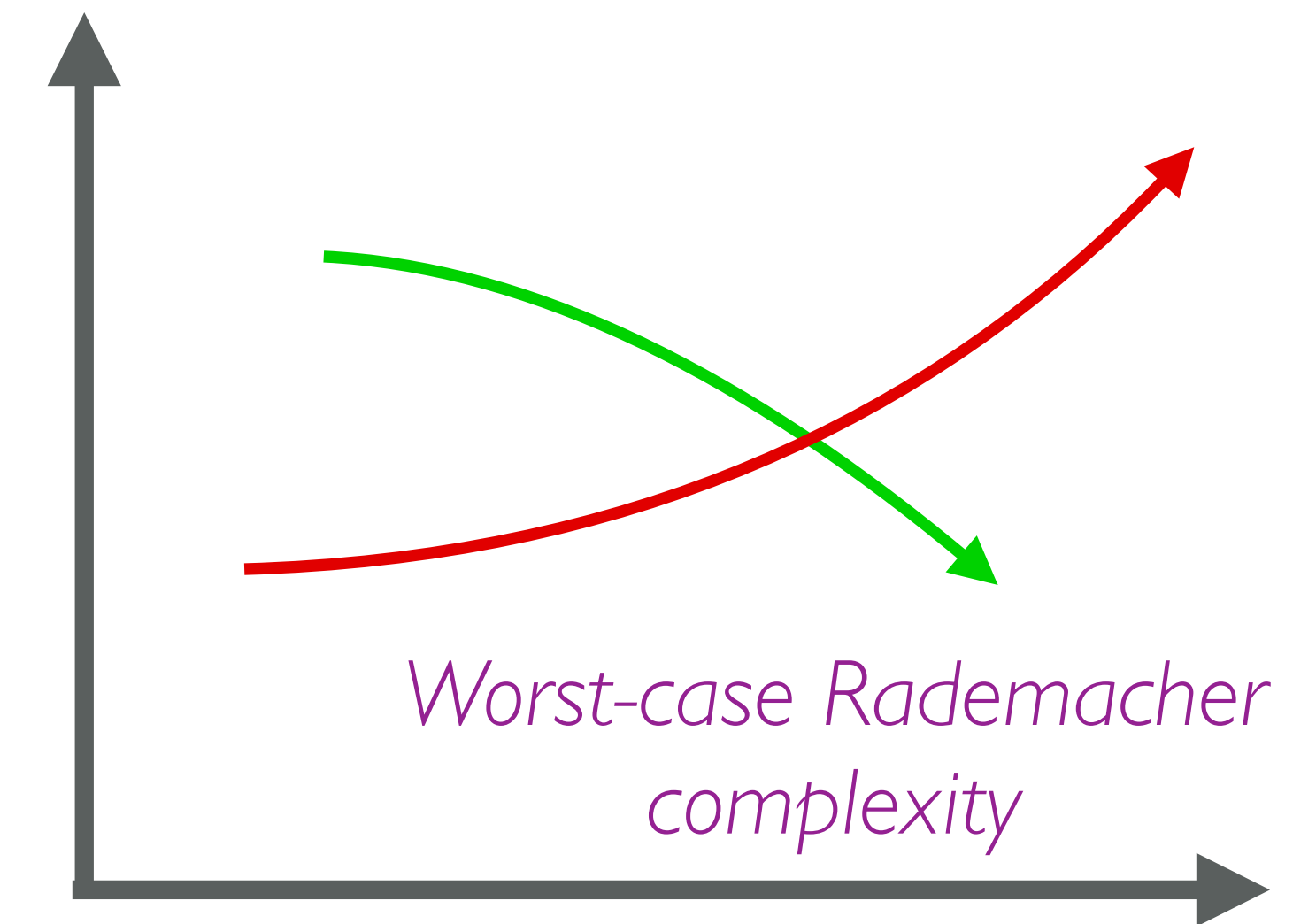
Approximation error + Estimation error



OUR LEARNING SETUP

Worst-case Rademacher complexity

$$\text{Rad}(\tilde{\mathcal{H}}) = \mathbb{E}_{\sigma_1, \dots, \sigma_m} \left[\sup_{\|\mathbf{w}\|_r \leq 1} \frac{1}{m} \sum_{i=1}^m \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \sigma_i \langle \mathbf{w}, \mathbf{x}'_i \rangle \right]$$



WORST-CASE RADEMACHER COMPLEXITY

Theorem [AFM 2020]

$$\text{Rad}(\tilde{\mathcal{H}}) \leq \text{Rad}(\mathcal{H}) + \frac{\epsilon \max(d^{1-\frac{1}{p}-\frac{1}{r}}, 1)}{2\sqrt{m}}$$

$$\text{Rad}(\tilde{\mathcal{H}}) \geq \max\left(\text{Rad}(\mathcal{H}), \frac{\epsilon \max(d^{1-\frac{1}{p}-\frac{1}{r}}, 1)}{2\sqrt{2m}}\right)$$

★ dimension dependent term, unless $1 - \frac{1}{p} - \frac{1}{r} \leq 0$

★ suggests *dual* regularization $r = \frac{p}{p-1}$

Setup: Hypothesis class

$\mathcal{H}_r = \{\mathbf{x} \rightarrow \langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{w}\|_r \leq 1\}$,
dataset $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$,
 ℓ_p perturbations of size ϵ .

worst-case Rademacher
complexity

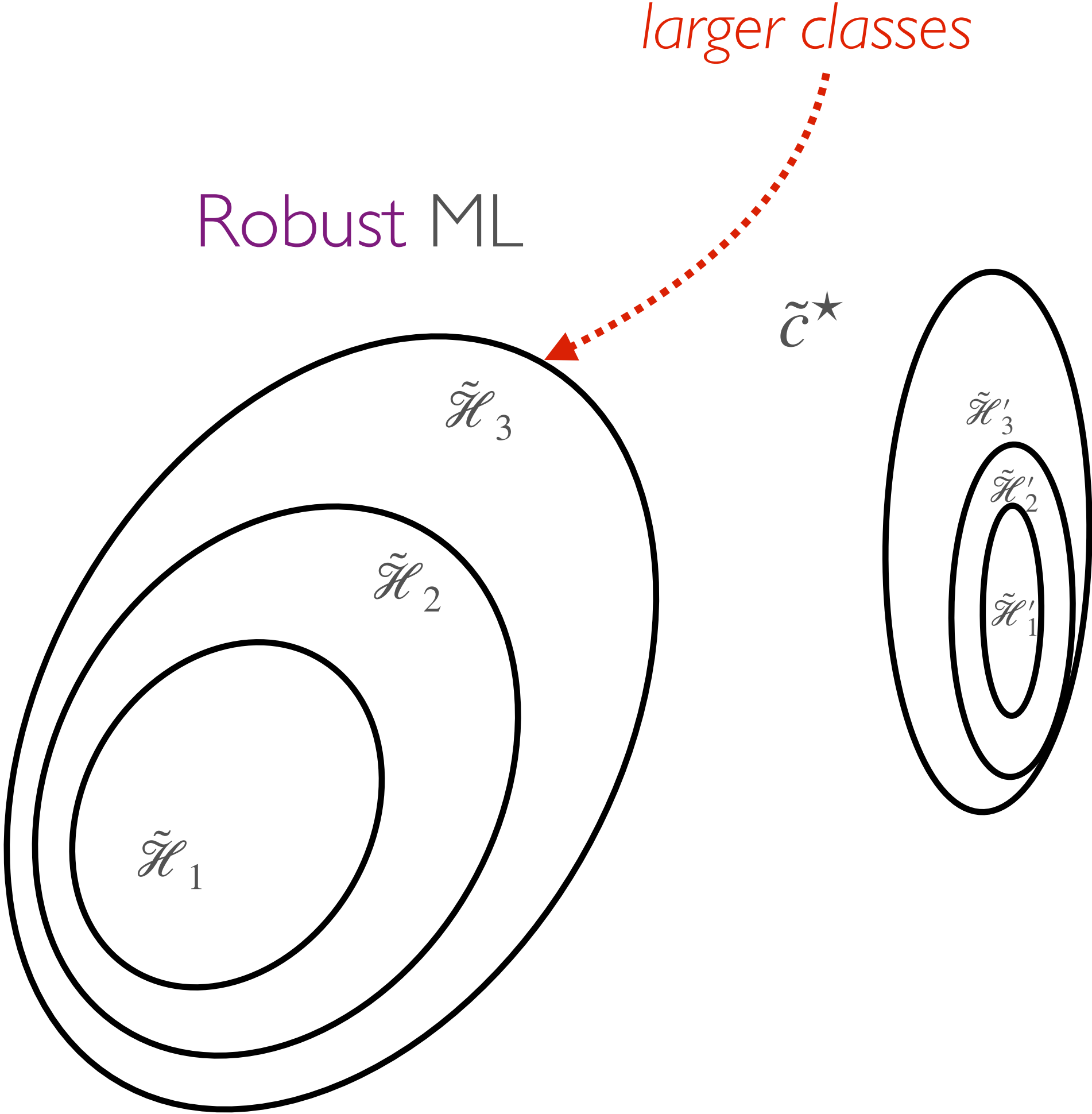
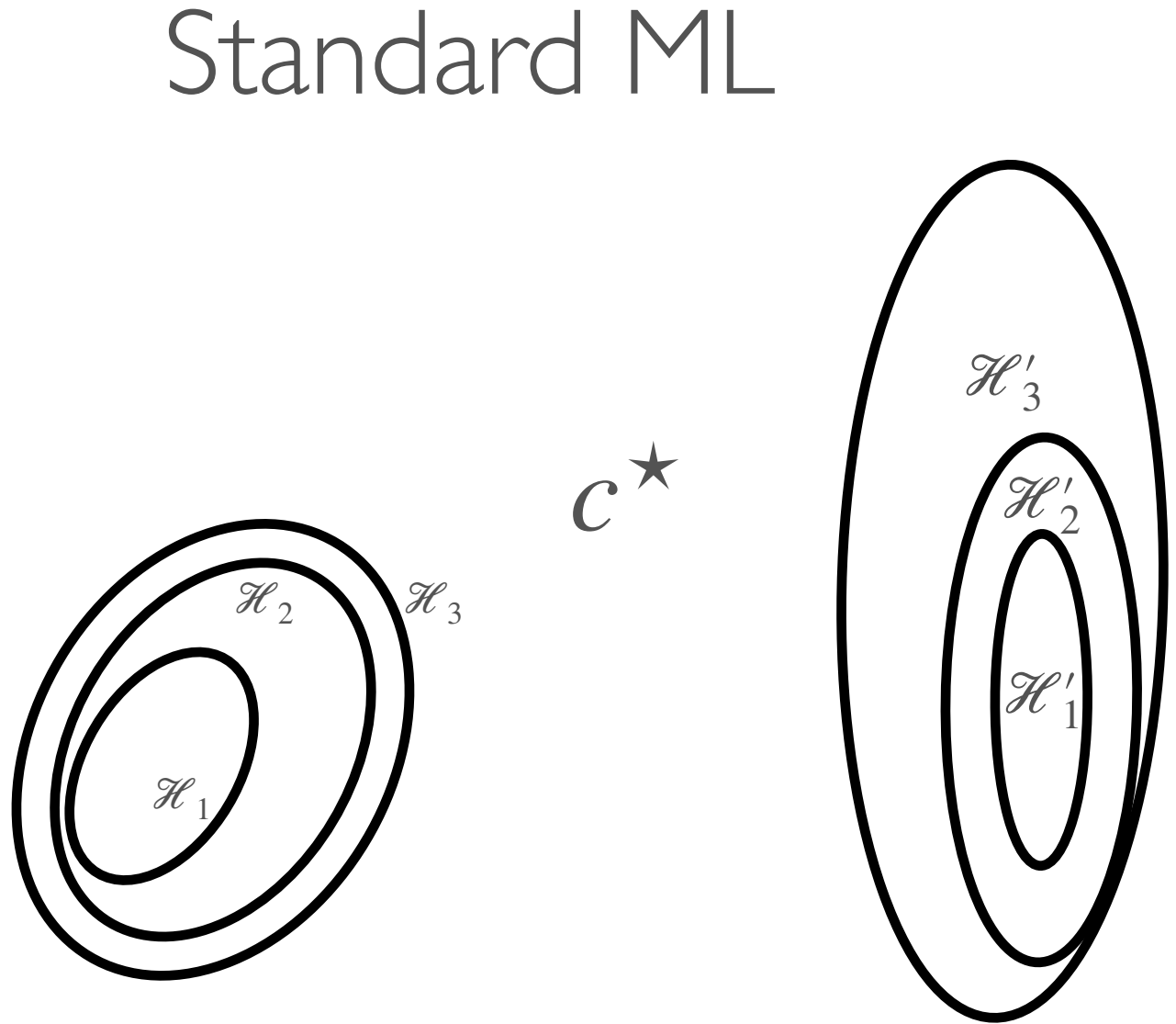
$$\text{Rad}(\tilde{\mathcal{H}}) = \mathbb{E}_{\sigma_1, \dots, \sigma_m} \left[\sup_{\|\mathbf{w}\|_r \leq 1} \frac{1}{m} \sum_{i=1}^m \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \sigma_i \langle \mathbf{w}, \mathbf{x}'_i \rangle \right]$$

$$\text{Rad}(\mathcal{H}) = \mathbb{E}_{\sigma_1, \dots, \sigma_m} \left[\sup_{\|\mathbf{w}\|_r \leq 1} \frac{1}{m} \sum_{i=1}^m \sigma_i \langle \mathbf{w}, \mathbf{x}_i \rangle \right]$$

Rademacher complexity

INTERIM SUMMARY

- Type of perturbation ℓ_p can influence “optimal” hypothesis class.
- Model selection becomes more important.



IMPLICIT BIAS OF OPTIMIZATION

In supervised learning, we do not get to choose regularization; the regularization chooses us (through optimization)

Optimization algorithm, architecture, loss function, step size, initialization implicitly perform model selection.

IMPLICIT BIAS IN ROBUST ERM: ALGORITHM (1/2)

Unregularized loss $\tilde{L}_S(\mathbf{w}) = \sum_{i=1}^m \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} e^{-y_i \langle \mathbf{w}, \mathbf{x}'_i \rangle}$. Algorithm: Steepest Descent w.r.t. $\|\cdot\|_r$ norm

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta \Delta \mathbf{w}_t,$$

where $\Delta \mathbf{w}_t \in \arg \min_{\|\mathbf{u}\|_r \leq \|\nabla L_S(\mathbf{w}_t)\|_{r^*}} \langle \mathbf{u}, \nabla L_S(\mathbf{w}_t) \rangle$.

Theorem

For any separable dataset and any initialization \mathbf{w}_0 , for $\eta \rightarrow 0^+$, it holds:

$$\lim_{t \rightarrow \infty} \min_{i \in [m]} \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \frac{y_i \langle \mathbf{w}_t, \mathbf{x}'_i \rangle}{\|\mathbf{w}_t\|_r} = \max_{\mathbf{w} \neq 0} \min_{i \in [m]} \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \frac{y_i \langle \mathbf{w}, \mathbf{x}'_i \rangle}{\|\mathbf{w}\|_r}.$$

IMPLICIT BIAS IN ROBUST ERM: ALGORITHM (2/2)

Theorem

For any separable dataset and any initialization

\mathbf{w}_0 , for $\eta \rightarrow 0^+$, it holds:

$$\lim_{t \rightarrow \infty} \min_{i \in [m]} \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \frac{y_i \langle \mathbf{w}_t, \mathbf{x}'_i \rangle}{\|\mathbf{w}_t\|_r} = \max_{\mathbf{w} \neq 0} \min_{i \in [m]} \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \frac{y_i \langle \mathbf{w}, \mathbf{x}'_i \rangle}{\|\mathbf{w}\|_r}.$$

Algorithm: Steepest Descent w.r.t. $\|\cdot\|_r$ norm

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta \Delta \mathbf{w}_t,$$

where $\Delta \mathbf{w}_t \in \arg \min_{\|\mathbf{u}\|_r \leq \|\nabla L_S(\mathbf{w}_t)\|_{r^*}} \langle \mathbf{u}, \nabla L_S(\mathbf{w}_t) \rangle$.

★ Equivalent to $\min_{\mathbf{w}} \|\mathbf{w}\|_r$ s.t. $\min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} y_i \langle \mathbf{w}, \mathbf{x}'_i \rangle \geq 1, \forall i \in [m]$.

★ For $r = 2$, robust ERM with GD $\rightarrow \|\mathbf{w}\|_2$ regularization (regardless of p).

★ Steepest descent for $r = p^* = \frac{p}{p-1} \rightarrow \|\mathbf{w}\|_{p^*}$ regularization.

IMPLICIT BIAS IN ROBUST ERM: ARCHITECTURE

Reparameterized linear model $f_{\text{diag}}(\mathbf{x}; \mathbf{w}) = \langle \mathbf{u}_+^2 - \mathbf{u}_-^2, \mathbf{x} \rangle$, $\mathbf{w} = [\mathbf{u}_+, \mathbf{u}_-]$.

$$\text{Unregularized loss } \tilde{L}_S(\mathbf{w}) = \sum_{i=1}^m \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} e^{-y_i f_{\text{diag}}(\mathbf{x}'_i; \mathbf{w})}.$$

Proposition (Informal)

Consider GD over $f_{\text{diag}}(\cdot; \mathbf{w})$ and assume $\tilde{L}_S(\mathbf{w}) \rightarrow 0$, then the implicit bias of \mathbf{w} corresponds to:

$$\min_{\mathbf{w}} \|\mathbf{w}\|_1 \text{ s.t. } \min_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} y_i \langle \mathbf{w}, \mathbf{x}'_i \rangle \geq 1, \forall i \in [m].$$



Same algorithm, different architecture \rightarrow different implicit bias.



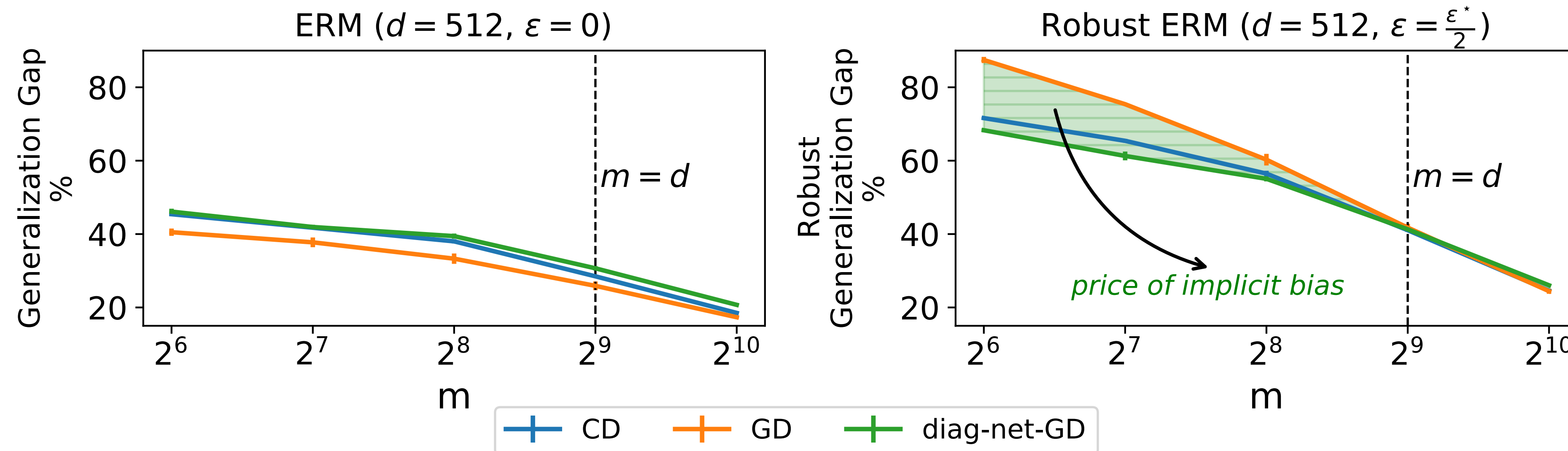
Favorable bias for $p = \infty$.

PRICE OF IMPLICIT BIAS

Coordinate Descent - CD:

Change only one coordinate per update (largest absolute gradient) - Steepest descent wrt ℓ_1 norm

- $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, I)$, $\mathbf{w}^* \sim \mathcal{N}(\mathbf{0}, I)$, $y = \text{sgn}(\langle \mathbf{w}^*, \mathbf{x} \rangle)$. Defend against ℓ_∞ perturbations.
- Find largest possible perturbation ϵ^* and use $\epsilon < \epsilon^*$.
- Compare GD on linear networks vs CD on linear networks vs GD on diagonal network.

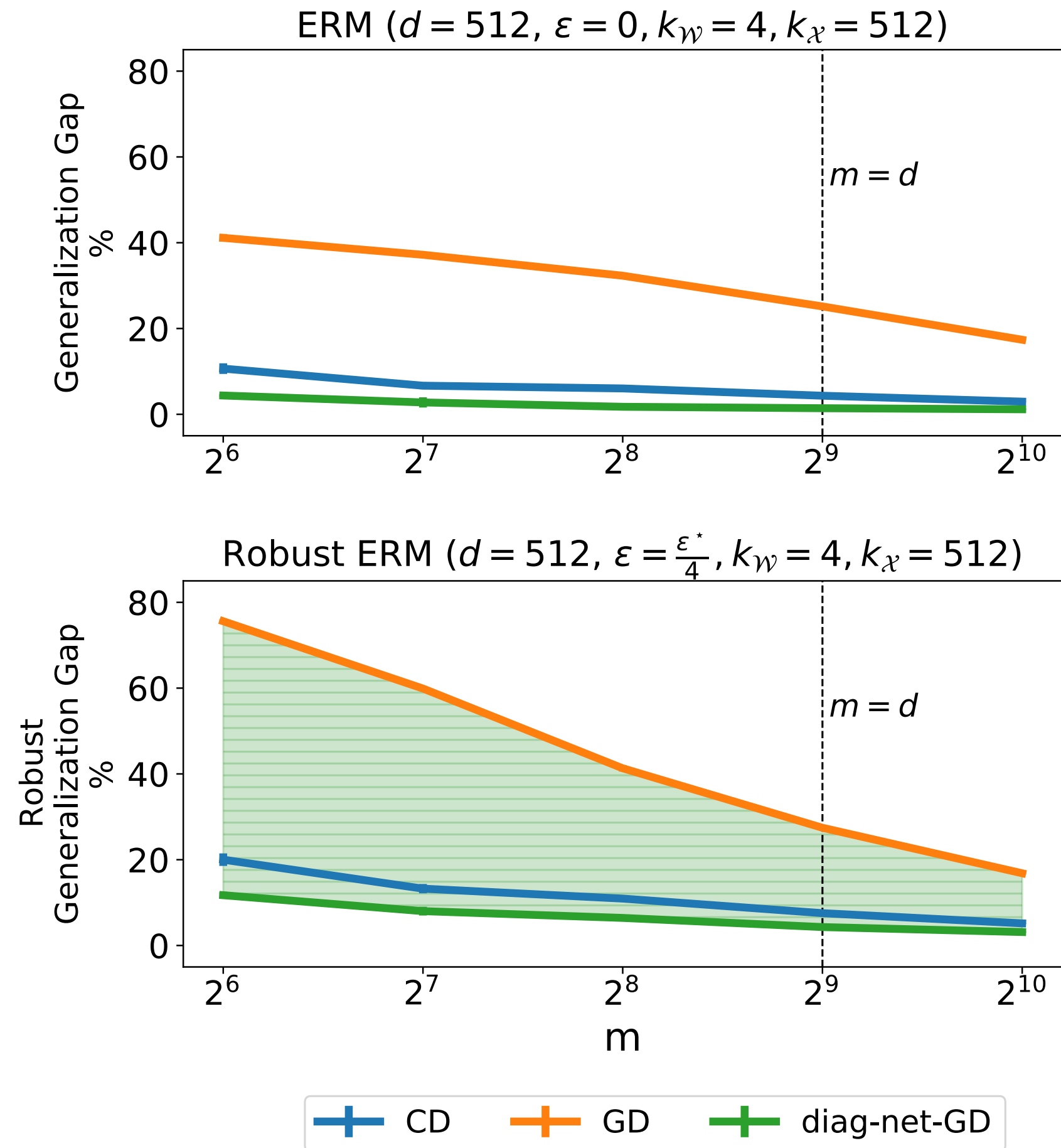


✓ Type of perturbation ℓ_p can influence “optimal” hypothesis class.

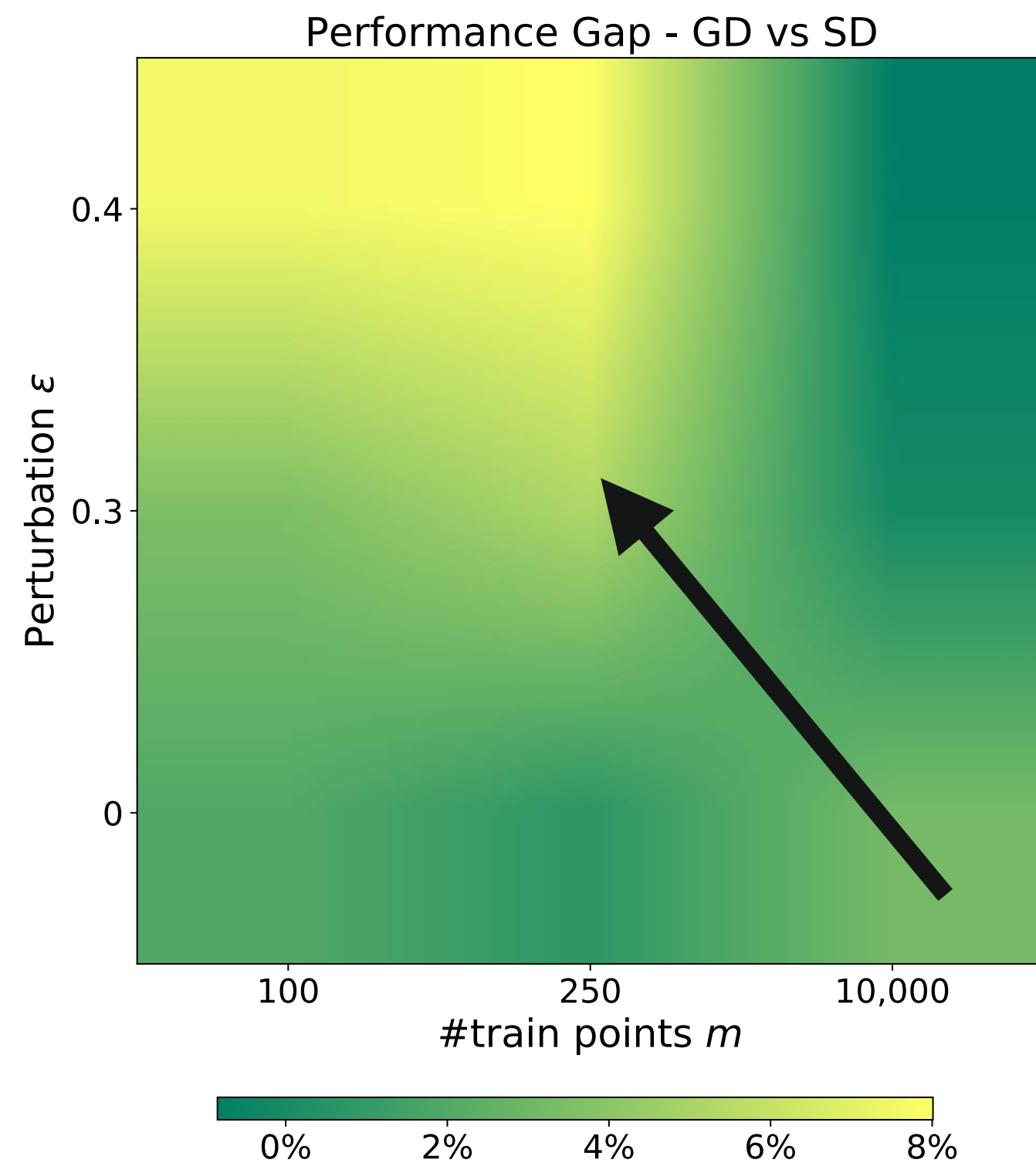
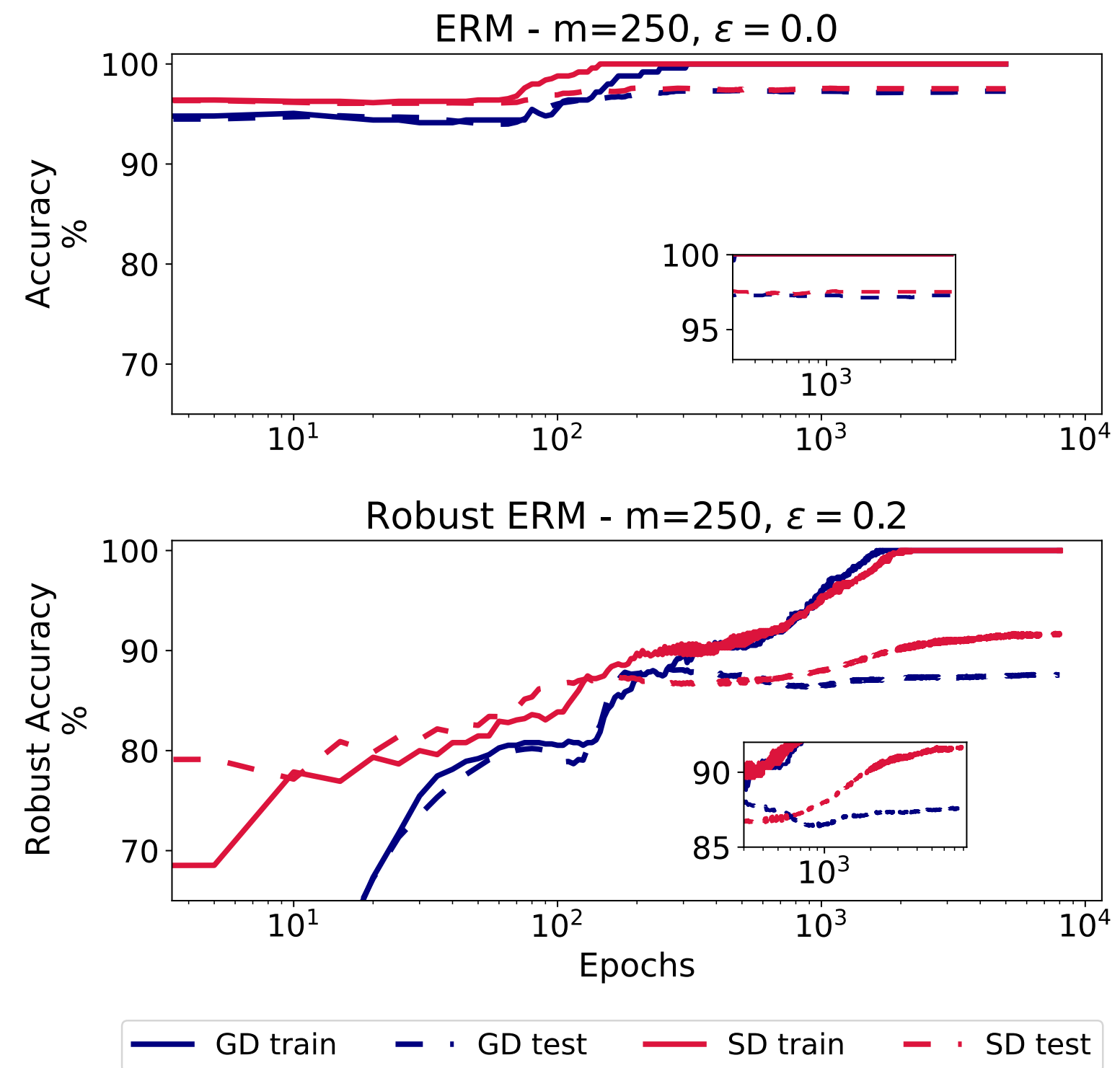
PRICE OF IMPLICIT BIAS

Same experiment with sparse teacher \mathbf{w}^\star

✓ *Model selection becomes more important.*



PRICE OF IMPLICIT BIAS IN NEURAL NETWORKS



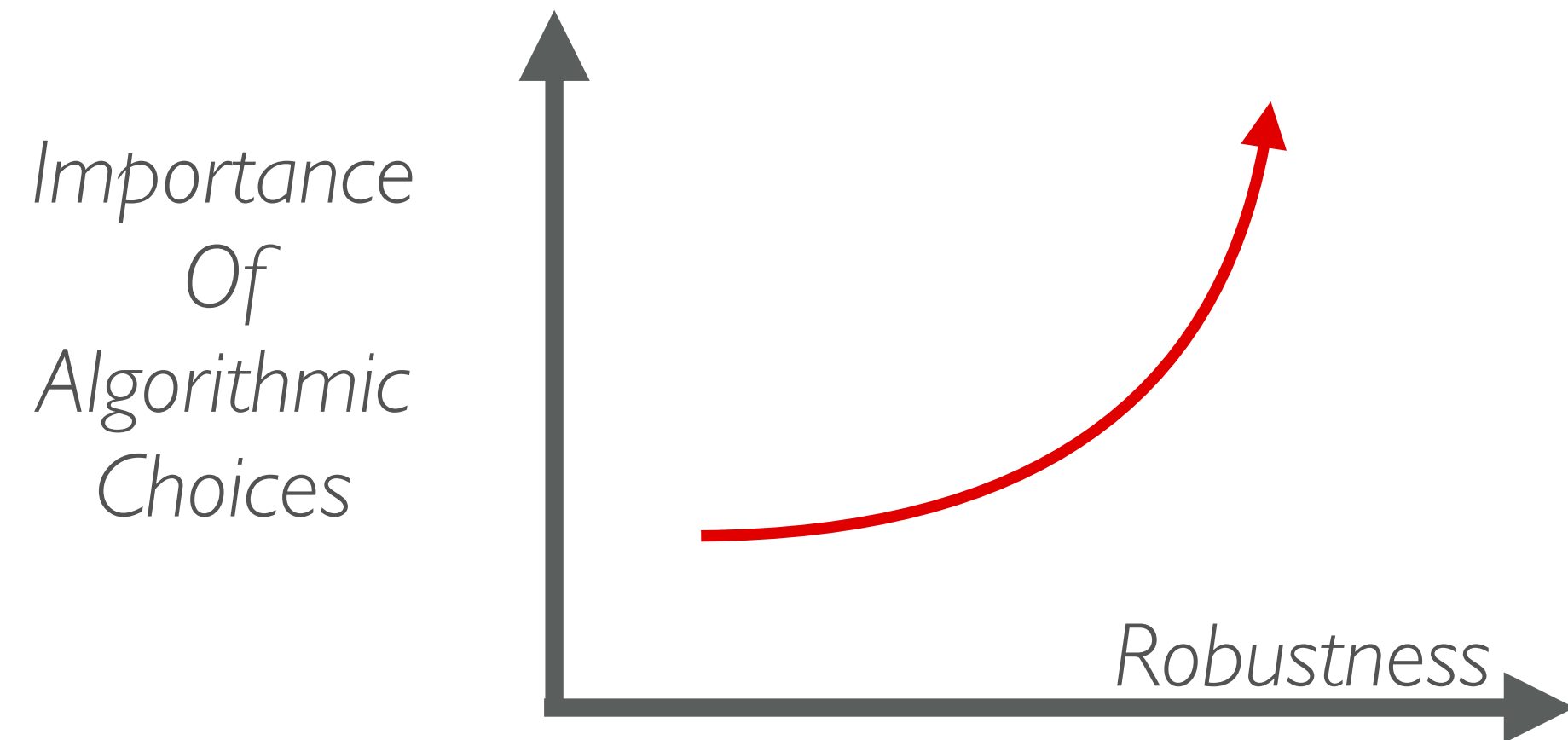
Sign Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \text{sgn}(\nabla L(\mathbf{w}_t))$$

[T., Vardi, Kempe. Flavors of Margin: Implicit Bias of Steepest Descent in Homogeneous Neural Networks]

[T., Frank, Srebro, Kempe. The Price of Implicit Bias in Adversarially Robust Generalization. NeurIPS 2024]

CONCLUSION



Robust training as
a testbed for
theories for deep
learning