

# Intro of Out-of-distribution Generalization

Dinghui Zhang

# What's "spurious" correlation?

## Common training examples

## Test examples

Waterbirds

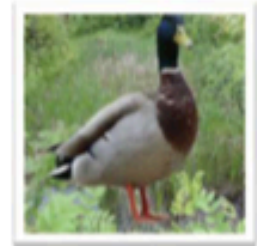
y: waterbird  
a: water  
background



y: landbird  
a: land  
background



y: waterbird  
a: land  
background

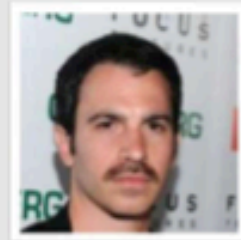


CelebA

y: blond hair  
a: female



y: dark hair  
a: male

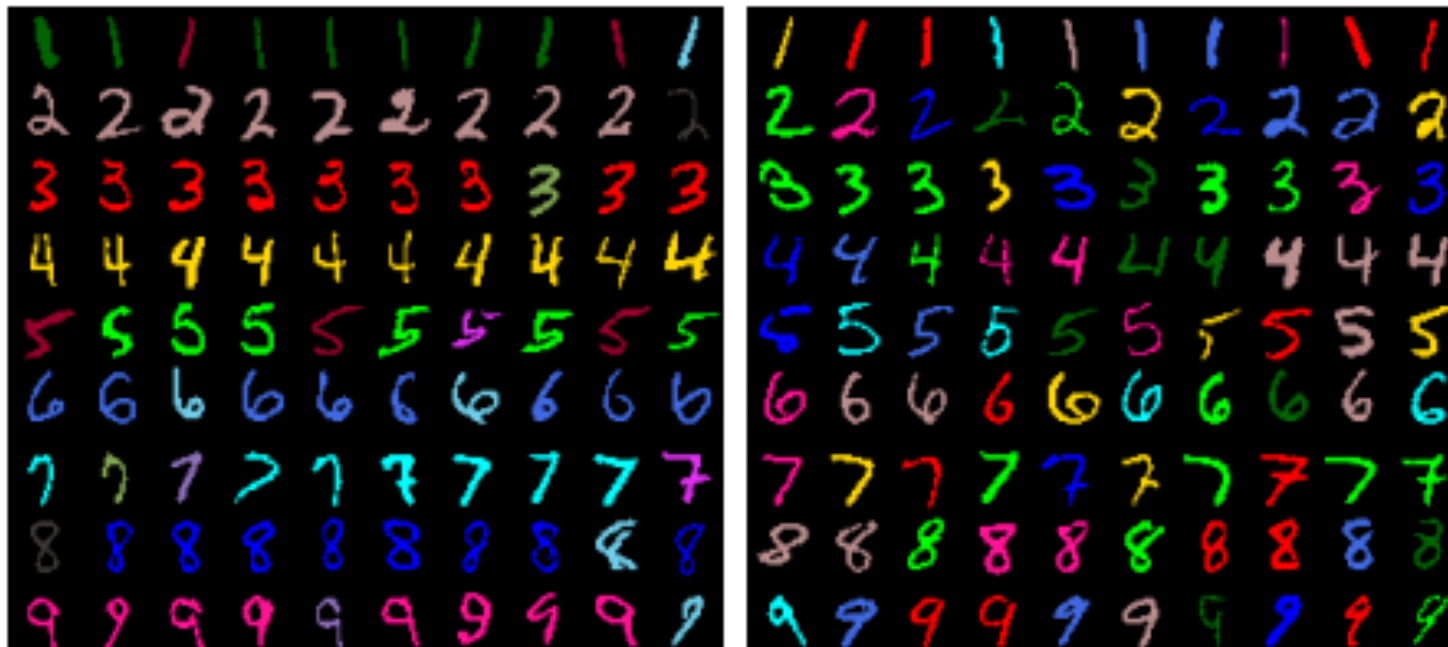


y: blond hair  
a: male



"true label" and "spurious label"

# What's "spurious" correlation?



“true label” and “spurious label”

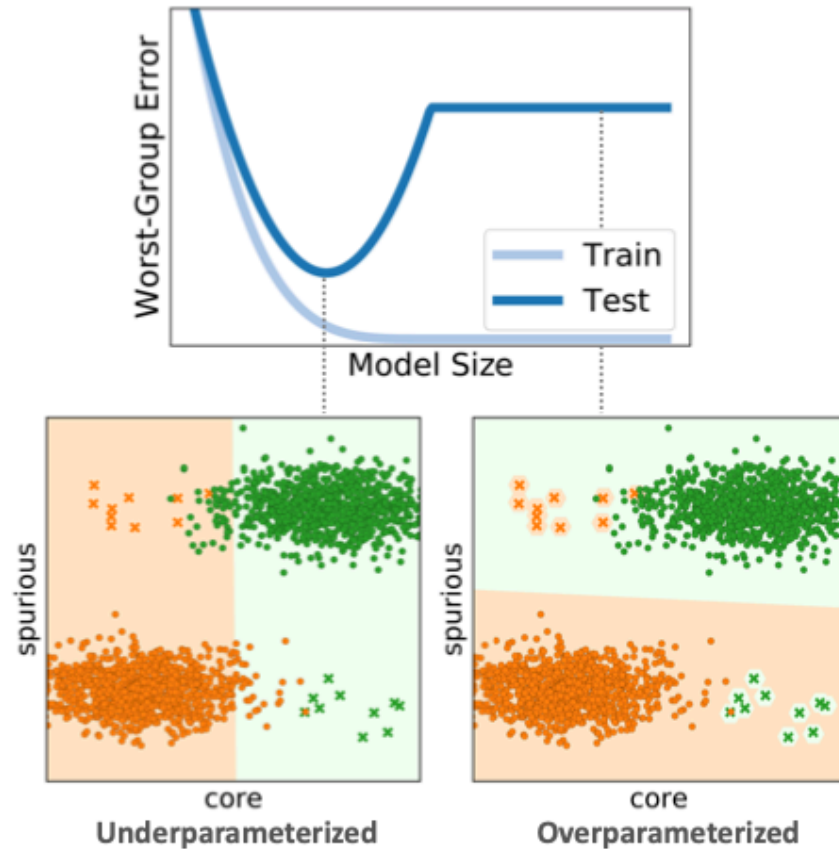
# GroupDRO

$$\hat{\theta}_{\text{DRO}} := \arg \min_{\theta \in \Theta} \left\{ \hat{\mathcal{R}}(\theta) := \max_{g \in \mathcal{G}} \mathbb{E}_{(x,y) \sim \hat{P}_g} [\ell(\theta; (x, y))] \right\},$$

DISTRIBUTIONALLY ROBUST NEURAL NETWORKS FOR GROUP SHIFTS: ON THE IMPORTANCE OF REGULARIZATION FOR WORST-CASE GENERALIZATION, Shiori Sagawa et al. ICLR2020

		Average Accuracy		Worst-Group Accuracy		
		ERM	DRO	ERM	DRO	
Standard Regularization	Waterbirds	Train	100.0	100.0	100.0	100.0
		Test	97.3	97.4	60.0	76.9
	CelebA	Train	100.0	100.0	99.9	100.0
		Test	94.8	94.7	41.1	41.1
	MultiNLI	Train	99.9	99.3	99.9	99.0
		Test	82.5	82.0	65.7	66.4
Strong $\ell_2$ Penalty	Waterbirds	Train	97.6	99.1	35.7	97.5
		Test	95.7	96.6	21.3	84.6
	CelebA	Train	95.7	95.0	40.4	93.4
		Test	95.8	93.5	37.8	86.7

# Overparameterization exacerbates spurious correlations

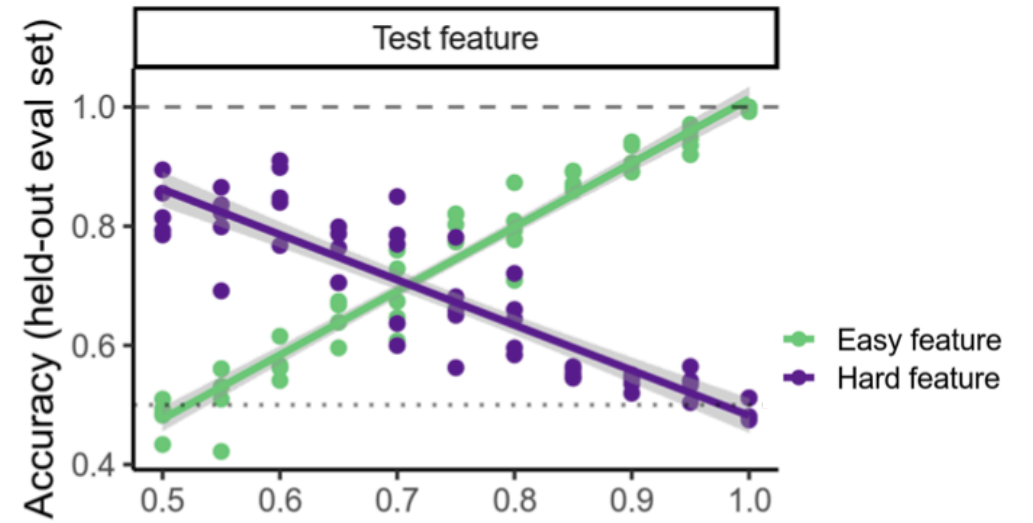
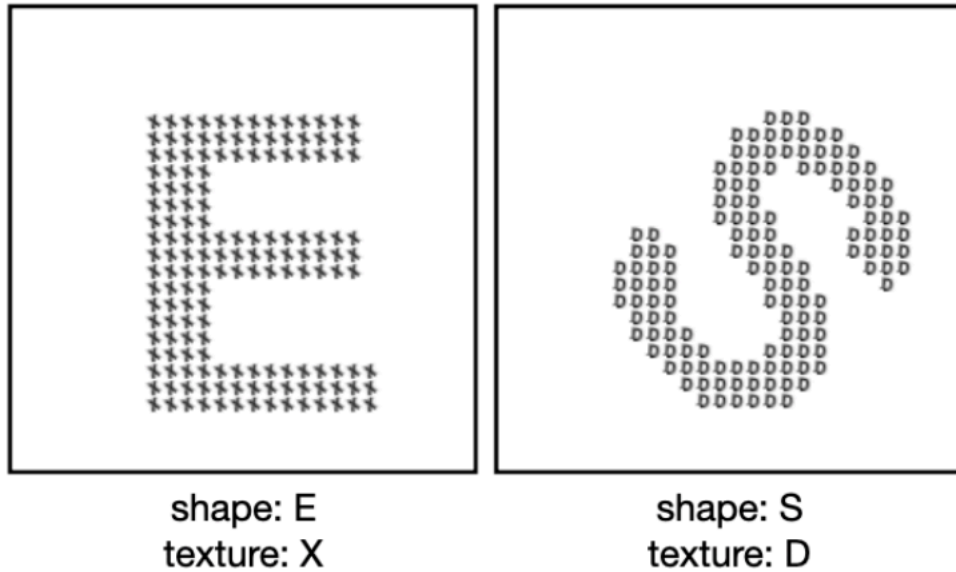


As model size grows, avg errors decrease, but worst group error increases

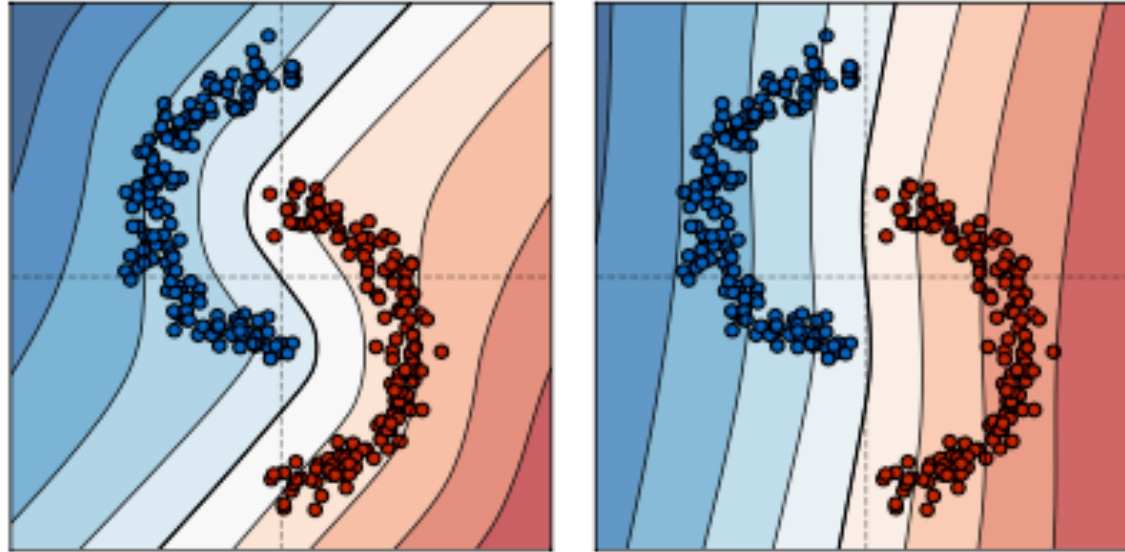
Reason: overparameterized models use spurious feature to classify

# Another name: Simplicity Bias

Navon



# Gradient Starvation

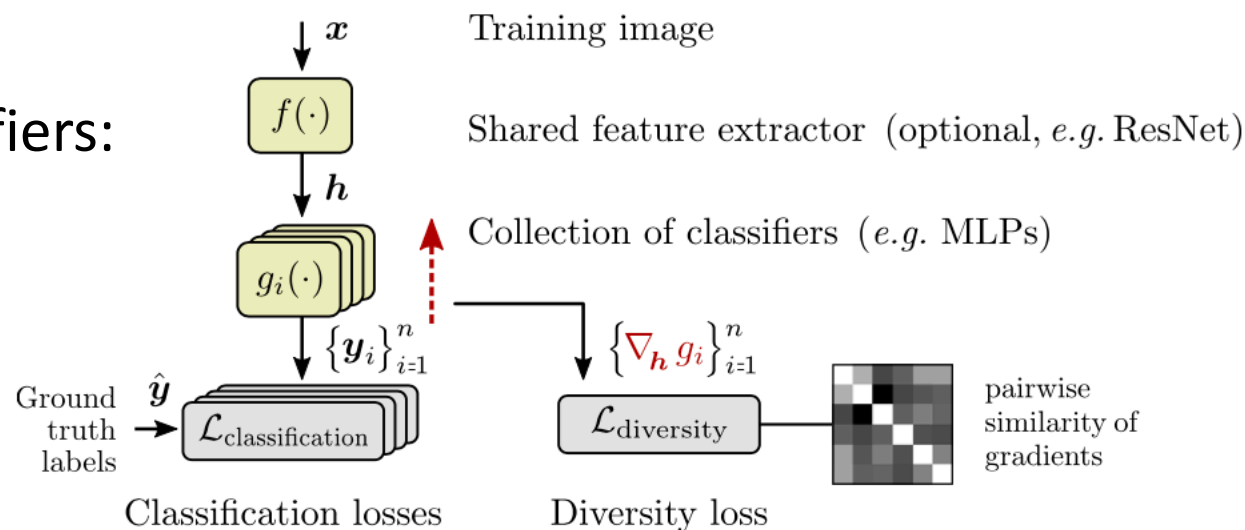


“overfitting” property of ERM



# Solution: Promote Diversity

Train multiple classifiers:



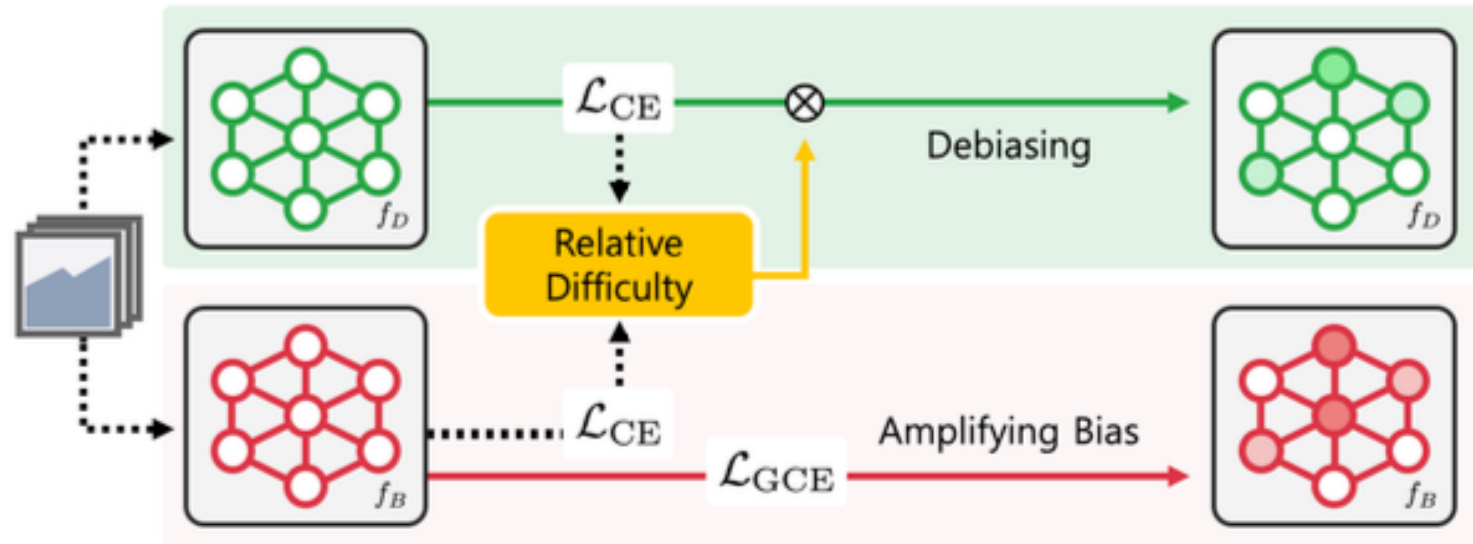
$$\min_{(\theta, \{\phi_i\})} \sum_i^n \mathcal{R}(g_{\phi_i} \circ f_{\theta}) + \lambda \sum_{i \neq j} \sum_k^K \delta_{g_{\phi_i}, g_{\phi_j}}(\mathbf{h}^k) \quad \delta_{g_{\phi_1}, g_{\phi_2}}(\mathbf{h}) = \nabla_{\mathbf{h}} g_{\phi_1}^*(\mathbf{h}) \cdot \nabla_{\mathbf{h}} g_{\phi_2}^*(\mathbf{h})$$

With post model selection method

Evading the Simplicity Bias: Training a Diverse Set of Models Discovers Solutions with Superior OOD Generalization

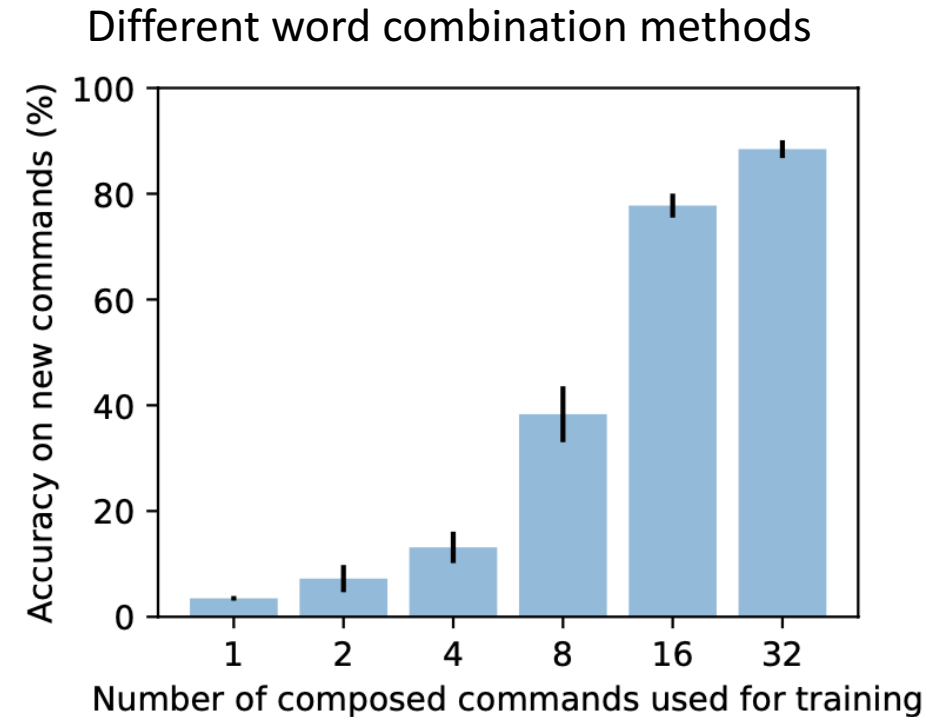
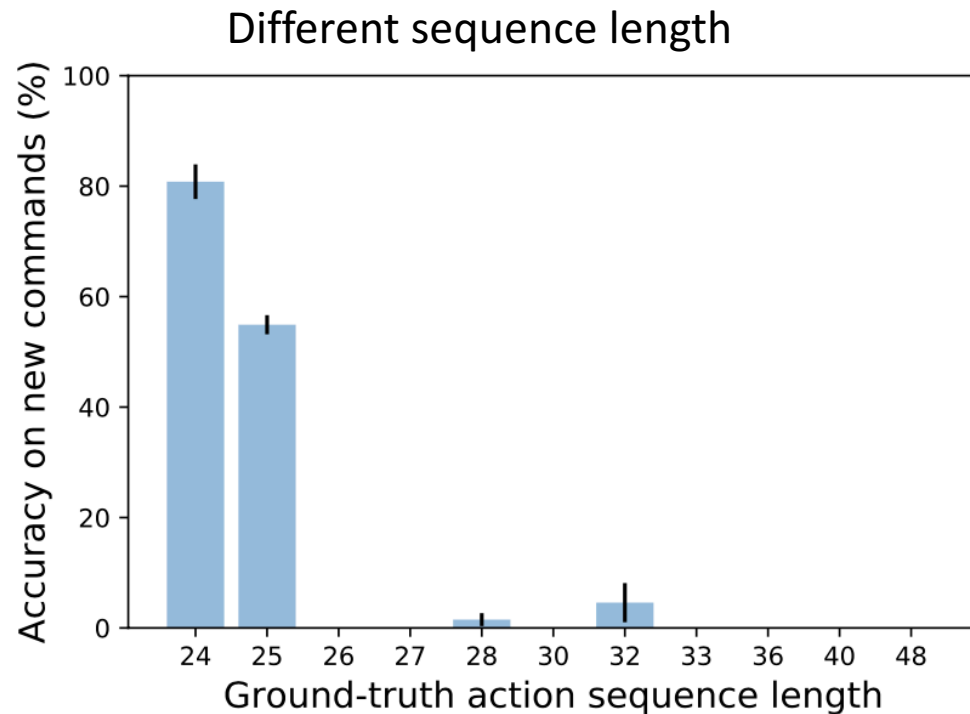
# Learning from Failure

- Setting: eg. No multiple domains
- 99% data: label & color has 1 to 1 corresponding
- 1% data: label & color has no corresponding



# NLP & CogSci: Compositionality

- How RNN generalize systematically under distribution shift



Generalization without Systematicity: On the Compositional Skills of Sequence-to-Sequence Recurrent Networks. ICML2018

# NLP: Debias

**MNLI synthetic:**

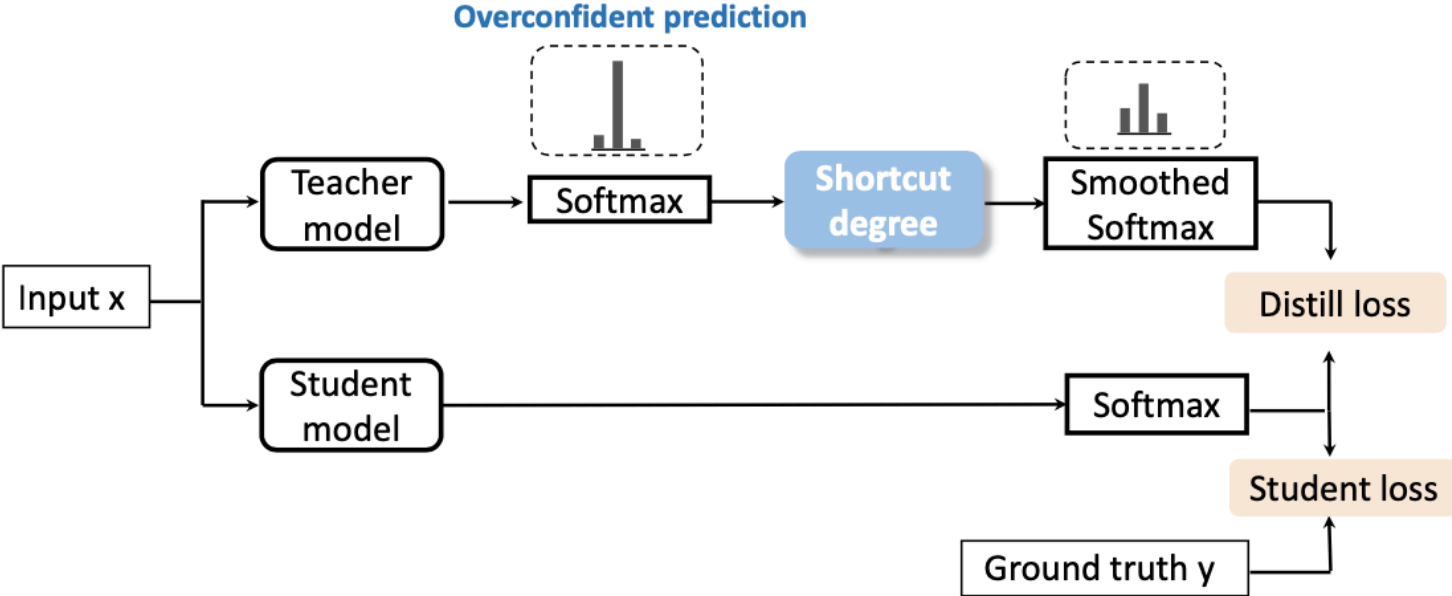
**premise:** What's truly striking, though, is that Jobs has never really let this idea go.

**orig. hypo.:** Jobs never held onto an idea for long.

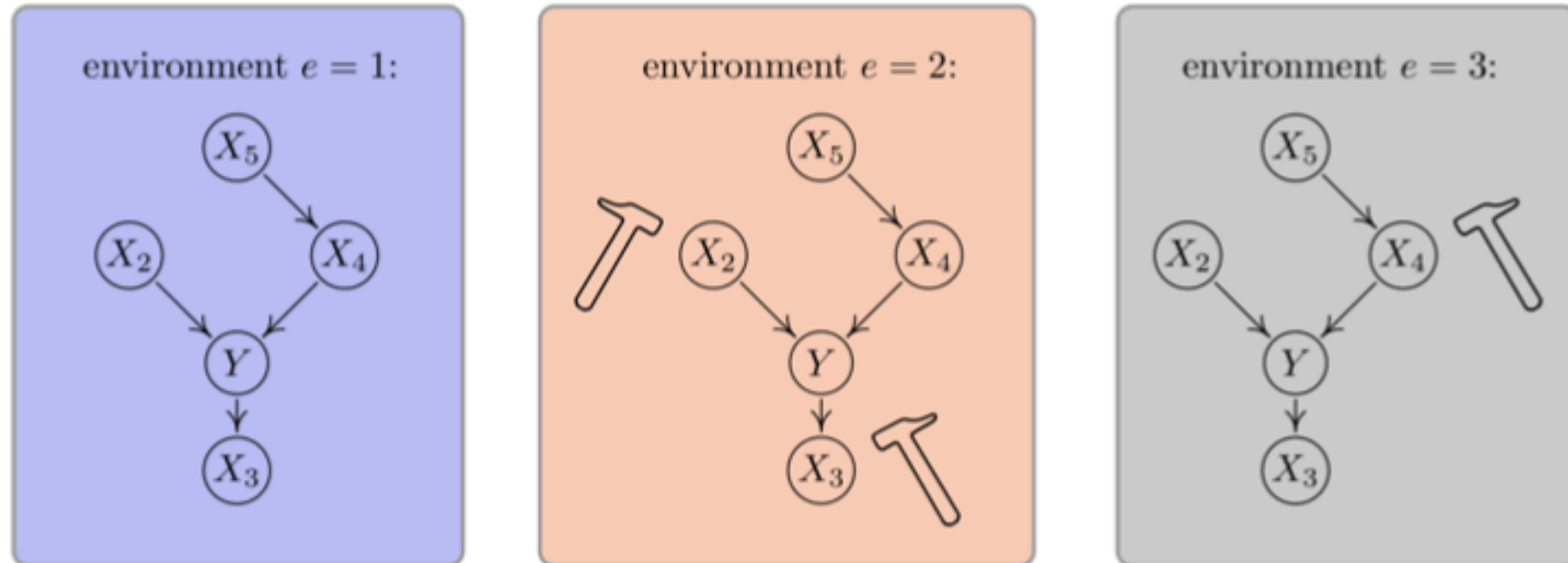
**biased:** 0 Jobs never held onto an idea for long.

**anti-biased:** 1 Jobs never held onto an idea for long.

**label:** 0 (contradiction)



# How to get rid of “spurious” feature? Or, how to do invariant learning



# Invariant Causal Prediction (ICP)

**Assumption 1 (Invariant prediction)** *There exists a vector of coefficients  $\gamma^* = (\gamma_1^*, \dots, \gamma_p^*)^t$  with support  $S^* := \{k : \gamma_k^* \neq 0\} \subseteq \{1, \dots, p\}$  that satisfies*

*for all  $e \in \mathcal{E}$  :  $X^e$  has an arbitrary distribution and*

$$Y^e = \mu + X^e \gamma^* + \varepsilon^e, \quad \varepsilon^e \sim F_\varepsilon \text{ and } \varepsilon^e \perp\!\!\!\perp X_{S^*}^e, \quad (3)$$

*where  $\mu \in \mathbb{R}$  is an intercept term,  $\varepsilon^e$  is random noise with mean zero, finite variance and the same distribution  $F_\varepsilon$  across all  $e \in \mathcal{E}$ .*

We will interchangeably use “domain” and “environment”.

# Causal Transfer Learning

## Algorithm 1: Subset search

**Inputs:** Sample  $(\mathbf{x}_i^k, y_i^k)_{i=1}^{n_k}$  for tasks  $k \in \{1, \dots, D\}$ , threshold  $\delta$  for independence test.

**Outputs:** Estimated invariant subset  $\hat{S}$ .

- 1 Set  $S_{acc} = \{\}$ ,  $MSE = \{\}$ .
- 2 **for**  $S \subseteq \{1, \dots, p\}$  **do**
- 3     linearly regress  $Y$  on  $\mathbf{X}_S$  and compute the residuals  $R_{\beta^{CS(S)}}$  on a validation set.
- 4     compute  $H = \text{HSIC}_b \left( (R_{\beta^{CS(S)}, i}, K_i)_{i=1}^n \right)$  and the corresponding p-value  $p^*$  (or the p-value from an alternative test, e.g., Levene test.).
- 5     **if**  $p^* > \delta$  **then**
- 6         compute  $\hat{\mathcal{E}}_{\mathbb{P}^{1, \dots, D}}(\beta^{CS(S)})$ , the empirical estimate of  $\mathcal{E}_{\mathbb{P}^{1, \dots, D}}(\beta^{CS(S)})$  on a validation set.
- 7          $S_{acc} \cdot \text{add}(S)$ ,  $MSE \cdot \text{add}(\hat{\mathcal{E}}_{\mathbb{P}^{1, \dots, D}}(\beta^{CS(S)}))$
- 8     **end**
- 9 **end**
- 10 Select  $\hat{S}$  according to *RULE*, see Section 3.4.

# Invariant Risk Minimization

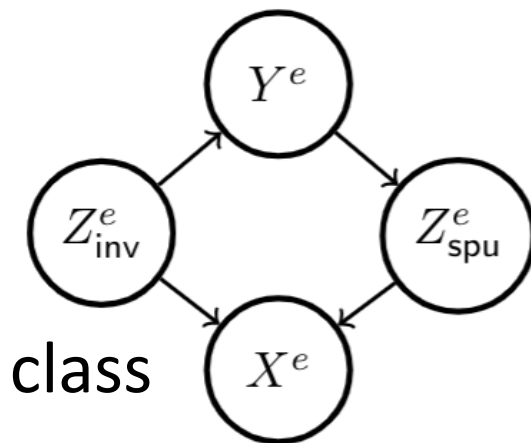
$$\begin{aligned} & \min_{\substack{\Phi: \mathcal{X} \rightarrow \mathcal{H} \\ w: \mathcal{H} \rightarrow \mathcal{Y}}} \sum_{e \in \mathcal{E}_{\text{tr}}} R^e(w \circ \Phi) \\ & \text{subject to } w \in \arg \min_{\bar{w}: \mathcal{H} \rightarrow \mathcal{Y}} R^e(\bar{w} \circ \Phi), \text{ for all } e \in \mathcal{E}_{\text{tr}}. \end{aligned} \quad (\text{IRM})$$

Require the classifier to be simultaneously optimal for all environments!

$$\min_{\Phi: \mathcal{X} \rightarrow \mathcal{Y}} \sum_{e \in \mathcal{E}_{\text{tr}}} R^e(\Phi) + \lambda \cdot \|\nabla_{w|w=1.0} R^e(w \cdot \Phi)\|^2, \quad (\text{IRMv1})$$



# ColoredMNIST



- Binary classification: 0~4 as positive class, 5~9 as negative class
- Each image is either red or green
- Domain1 (train): In all positive images, 70% are red; in all negative images, 30% are red.

Algorithm	Acc. train envs.	Acc. test env.
ERM	$87.4 \pm 0.2$	$17.1 \pm 0.6$
<b>IRM (ours)</b>	$70.8 \pm 0.9$	<b><math>66.9 \pm 2.5</math></b>
Random guessing (hypothetical)	50	50
Optimal invariant model (hypothetical)	75	75
ERM, grayscale model (oracle)	$73.5 \pm 0.2$	$73.0 \pm 0.4$

# Game theory formulation

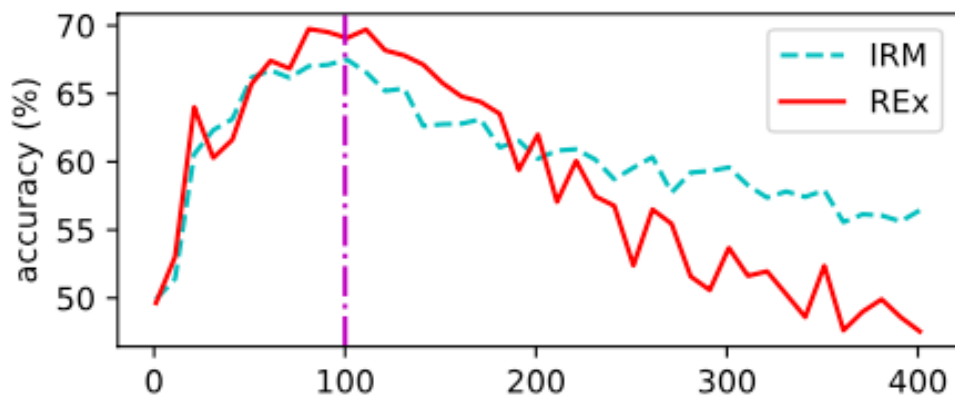
$$\begin{aligned} & \min_{\Phi \in \mathcal{H}_\Phi, w^{av} \in \mathcal{H}_w} \sum_{e \in \mathcal{E}_{tr}} R^e(w^{av} \circ \Phi) \\ \text{s.t. } & w^e \in \arg \min_{\bar{w}^e \in \mathcal{H}_w} R^e \left( \frac{1}{|\mathcal{E}_{tr}|} \left[ \bar{w}^e + \sum_{q \neq e} w^q \right] \circ \Phi \right), \forall e \in \mathcal{E}_{tr} \end{aligned}$$

A game between many classifiers

$$\begin{aligned} & \min_{\Phi \in \mathcal{H}_\Phi, w^{av} \in \mathcal{H}_w} \sum_{e \in \mathcal{E}_{tr}} R^e(w^{av} \circ \Phi) \\ \text{s.t. } & R^e \left( \frac{1}{|\mathcal{E}_{tr}|} \left[ w^e + \sum_{q \neq e} w^q \right] \circ \Phi \right) \\ & \leq R^e \left( \frac{1}{|\mathcal{E}_{tr}|} \left[ \bar{w}^e + \sum_{q \neq e} w^q \right] \circ \Phi \right) \quad \forall \bar{w}^e \in \mathcal{H}_w \quad \forall e \in \mathcal{E}_{tr} \end{aligned} \tag{3}$$

# REx

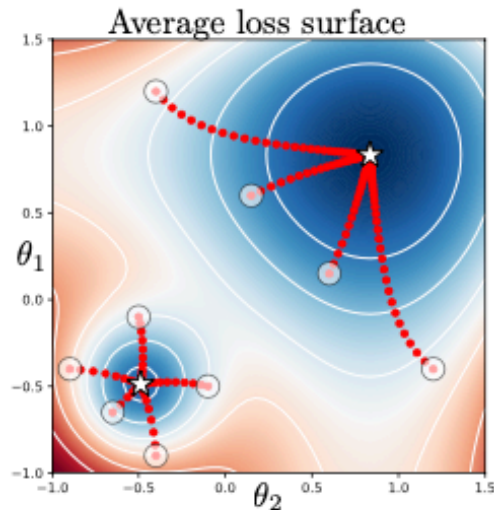
$$\mathcal{R}_{V\text{-REx}} \doteq \beta \text{Var}(\{\mathcal{R}_1, \dots, \mathcal{R}_m\}) + \sum_{e=1}^m \mathcal{R}_e$$



Feature selection effects

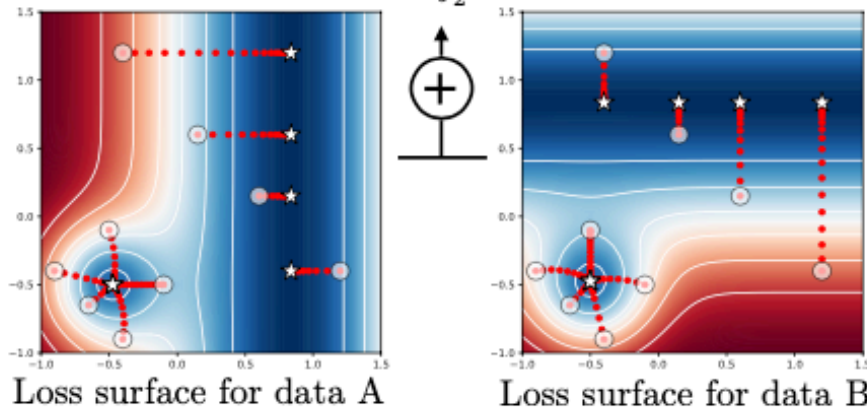
**Theorem 1.** *Given a Linear SEM,  $X_i \leftarrow \sum_{j \neq i} \beta_{(i,j)} X_j + \varepsilon_i$ , with  $Y \doteq X_0$ , and a predictor  $f_\beta(X) \doteq \sum_{j:j>0} \beta_j X_j + \varepsilon_j$  that satisfies REx (with mean-squared error) over a perturbation set of domains that contains 3 distinct  $do()$  interventions for each  $X_i : i > 0$ . Then  $\beta_j = \beta_{0,j}, \forall j$ .*

# Learning explanations that are hard to vary



$$\mathcal{C}^\epsilon(\theta^*) := - \max_{(e, e') \in \mathcal{E}^2} \max_{\theta \in N_{e, \theta^*}^\epsilon} |\mathcal{L}_{e'}(\theta) - \mathcal{L}_e(\theta)|.$$

Find the solution where the local geometry is invariant  
(2 order information)



“and mask”

a *threshold*  $\tau \in [0, 1]$

$$[\tilde{m}_\tau]_j = \mathbb{1} [\tau d \leq |\sum_e \text{sign}([\nabla \mathcal{L}_e]_j)|].$$

$$m_t(\theta) \odot \nabla \mathcal{L}(\theta)$$

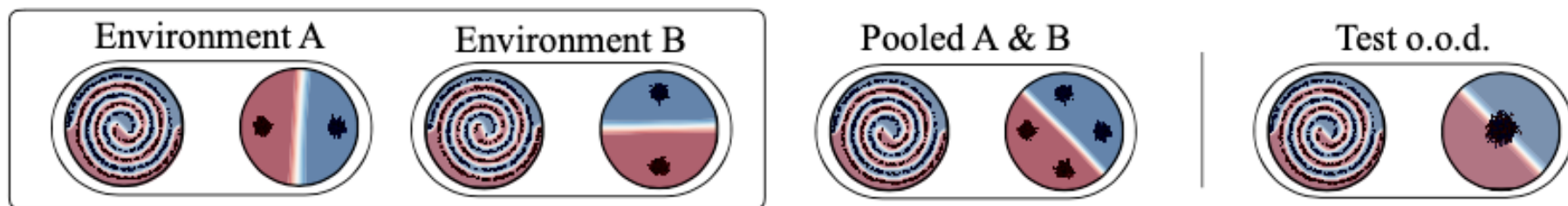


Figure 5: A 4-dimensional instantiation of the synthetic memorization dataset for visualization. Every example is a dot in both circles, and it can be classified by finding either of the “oracle” decision boundaries shown.

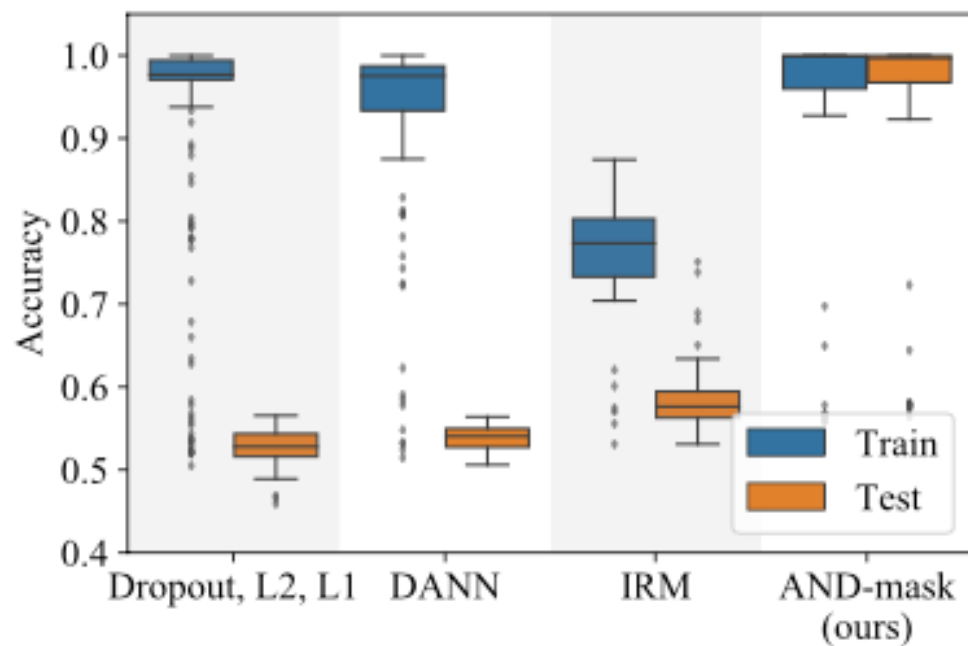
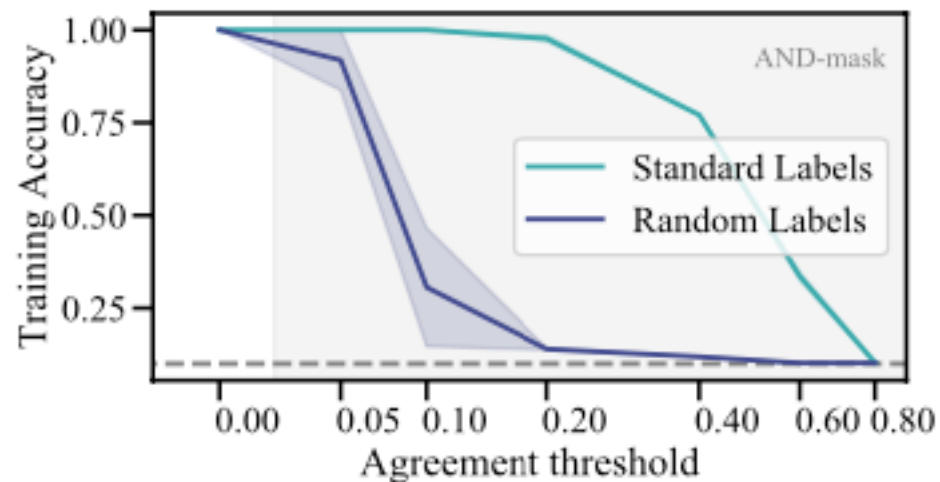


Figure 6: Results on the synthetic dataset.

# CIFAR10 random label



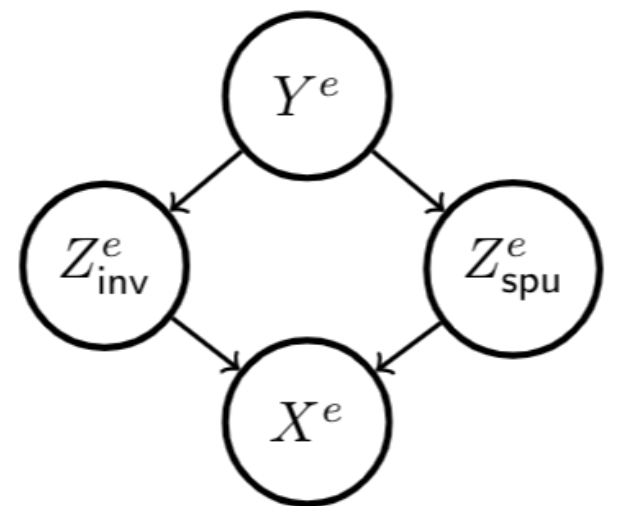
**Figure 8:** As the AND-mask threshold increases, memorization on CIFAR-10 with random labels is quickly hindered.

# Risks of IRM

$$y = \begin{cases} 1, & \text{w.p. } \eta \\ -1, & \text{otherwise.} \end{cases}$$

$$z_c \sim \mathcal{N}(y \cdot \mu_c, \sigma_c^2 I), \quad z_e \sim \mathcal{N}(y \cdot \mu_e, \sigma_e^2 I)$$

$$x = f(z_c, z_e).$$



$$\mathcal{R}^e(\Phi, \hat{\beta}) := \mathbb{E}_{(x,y) \sim p^e} \left[ \ell(\sigma(\hat{\beta}^T \Phi(x)), y) \right]$$

$$\min_{\Phi, \hat{\beta}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \left[ \mathcal{R}^e(\Phi, \hat{\beta}) + \lambda \|\nabla_{\hat{\beta}} \mathcal{R}^e(\Phi, \hat{\beta})\|_2^2 \right]$$



**Proposition 4.1.** *Suppose the observed data are generated according to Equations 1-3. Then recovering the (parametrized) invariant classifier  $\Phi(x) = [z_c]$  and  $\hat{\beta} = [\beta_c, \beta_0]$  is a stationary point for Equation 4.*

**Theorem 5.1** (Linear case). *Assume  $f$  is linear. Suppose we observe  $E$  training environments. Then the following hold:*

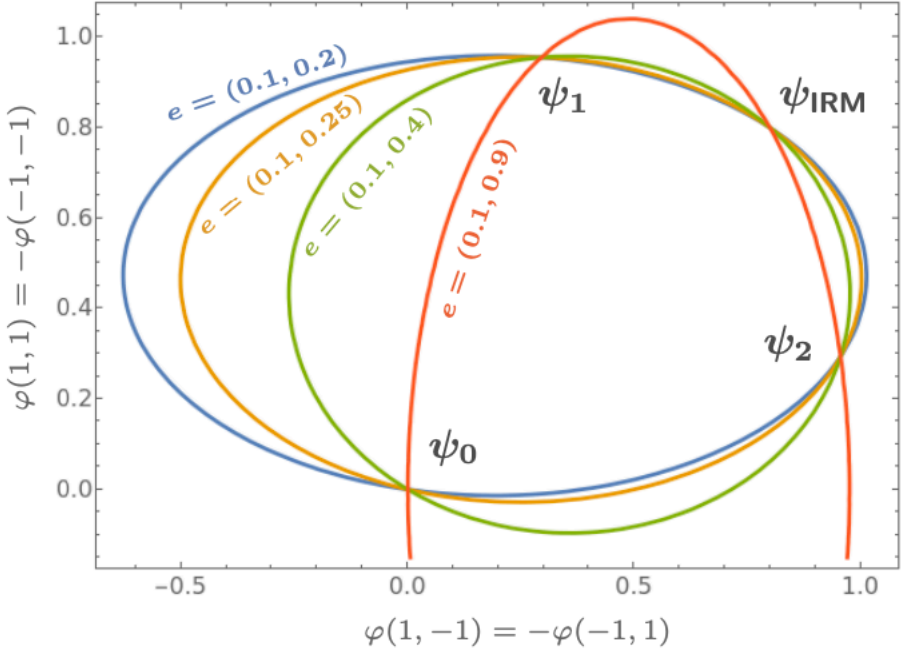
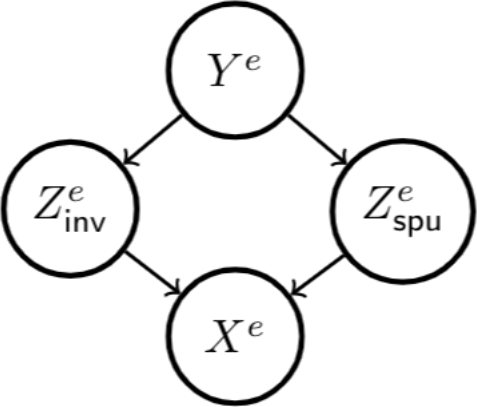
- 1. Suppose  $E > d_e$ . Under mild non-degeneracy conditions, any linear featurizer  $\Phi$  with an invariant optimal regression vector  $\hat{\beta}$  uses only invariant features, and it therefore has identical risk on all possible environments.*
- 2. If  $E \leq d_e$  and the environmental means  $\mu_e$  are linearly independent, then there exists a linear  $\Phi$  with  $\text{rank}(\Phi) = d_c + d_e + 1 - E$  whose output depends on the environmental features, yet the optimal classifier on top of  $\Phi$  is invariant. Further, both the logistic and 0-1 risks of this  $\Phi$  and its corresponding  $\hat{\beta}$  are strictly lower than those of the invariant classifier.*

# Equivalence to fairness

- Env index can be seen as sensitive attributes  $\mathbb{E}[y|\Phi(x)] = h, e]$
- Settings without domain label:
  1. Input *reference model*  $\tilde{\Phi}$ ;
  2. Fix  $\Phi \leftarrow \tilde{\Phi}$  and fully optimize the inner loop of (EIL) to infer environments  $\tilde{\mathbf{q}}_i(e) = \tilde{q}(e|x_i, y_i)$ ;
  3. Fix  $\mathbf{q} \leftarrow \tilde{\mathbf{q}}$  and fully optimize the outer loop to yield the new model  $\Phi$ .

# Does IRM Capture Invariance?

$$\begin{aligned}
 Y &\leftarrow \text{Rad}(0.5), \\
 X_1 &\leftarrow Y \cdot \text{Rad}(\alpha_e), \\
 X_2 &\leftarrow Y \cdot \text{Rad}(\beta_e),
 \end{aligned}
 \quad (\text{Two-Bit-Envs})$$

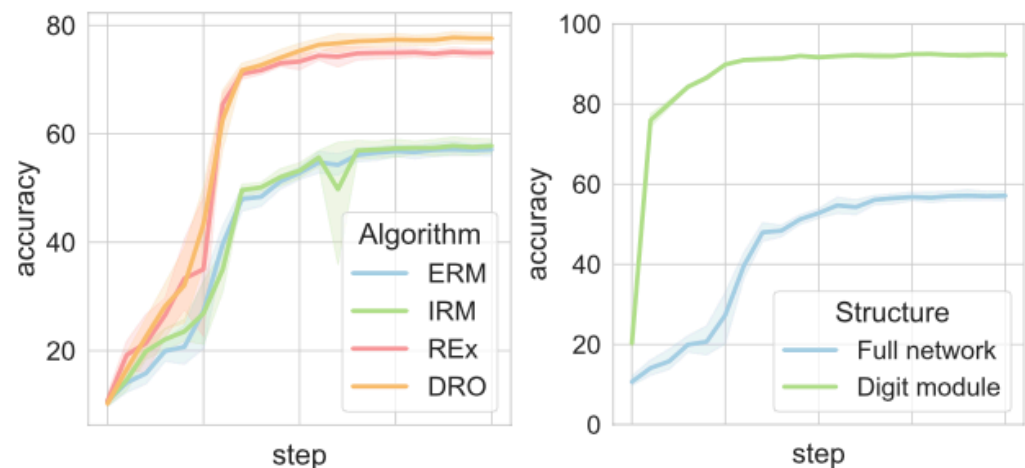


**Observation 2.** Under Setting A, a representation  $\varphi : \mathcal{X} \rightarrow \mathcal{Z}$  is invariant over  $\mathcal{E}$  if and only if for all  $e_1, e_2 \in \mathcal{E}$ , it holds that

$$\mathbb{E}_{\mathcal{D}_{e_1}}[Y \mid \varphi(X) = z] = \mathbb{E}_{\mathcal{D}_{e_2}}[Y \mid \varphi(X) = z]$$

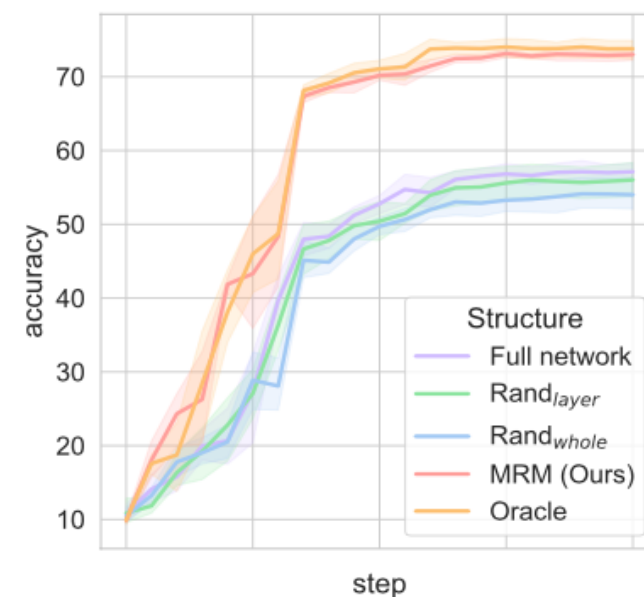
# Invariant subnetwork property

Invariant subnetwork exists in normally trained large network:

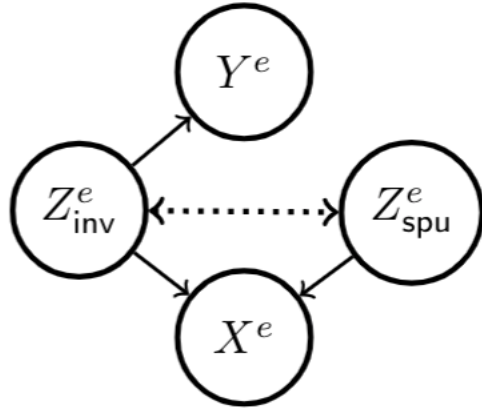


(a) Accuracy of baselines. (b) Accuracy of module in ERM.

Design algorithm to utilize this property:



# Information bottleneck (IB) principle



$$\min_{w \in \mathbb{R}^{k \times r}, \Phi \in \mathbb{R}^{r \times d}} \frac{1}{|\mathcal{E}_{tr}|} \sum_{e \in \mathcal{E}_{tr}} h^e(w \cdot \Phi)$$

$$\text{s.t. } \frac{1}{|\mathcal{E}_{tr}|} \sum_{e \in \mathcal{E}_{tr}} R^e(w \cdot \Phi) \leq r^{\text{th}}, \quad w \in \arg \min_{\tilde{w} \in \mathcal{H}_w} R^e(\tilde{w} \cdot \Phi), \forall e \in \mathcal{E}_{tr}$$

## Theorem 4. *IB-IRM and IB-ERM vs IRM and ERM*

- **Fully informative invariant features.** Suppose that the data  $\forall e \in \mathcal{E}_{all}$  follows Assumption 2. Assume that the invariant features are separable, bounded, and satisfy support overlap (Assumptions 3,5 and 7 hold). Also,  $\forall e \in \mathcal{E}_{tr} Z_{\text{spu}}^e \leftarrow AZ_{\text{inv}}^e + W^e$ , where  $W^e$  is continuous, bounded, zero mean noise. Every solution of IB-IRM (equation (6),  $\ell$  is 0-1 loss,  $r^{\text{th}} = q$ ), and IB-ERM solves OOD generalization (equation (1)) but ERM and IRM fail.

- More papers at [https://sites.google.com/site/irinarish/ood\\_generalization](https://sites.google.com/site/irinarish/ood_generalization)
- Thank you very much!