

Week - 4

Domain Generalization

- Sameer Jain
 - Pranav Mani

Causal Inference using Invariant Prediction Identification and Confidence intervals

Causal Model \rightarrow Invariant conditional distributions

Goal: Invariant conditional distributions \rightarrow Causal Model

* E : set of environments. $e \in E$ specifies a distribution over x, y . Alternate denote x, y by x^e, y^e

* The distributions are valid interventions on the causal structure.

↓
What's fair game?
• Mess with anything
except SEM(Y)

* $S \subseteq \{1 \dots p\}$, $p = |\mathcal{X}|$, $X_S : \{x_i | i \in S\}$

* $S^* : \{i | x_i \in \text{Pa}(Y)\}$

* for allowed interventions: $P(Y | \text{Pa}(Y))$ is an invariant structure across E

for linear SEM

ASSUMPTION - I :

$$\exists Y^* \text{ with support } S^* \text{ s.t. } \forall e \in E \\ y = x_{y^*}^e + \epsilon^e. \quad \epsilon^e \sim F_e, \quad \epsilon^e \perp\!\!\!\perp X_S^*$$

- * By definition assumption-I is satisfied by $P_{\alpha}(Y)$
- * But many S can satisfy assumption
reg) X, Y : jointly Gaussian
- * Frame as hypothesis

$$H_{0,r,S}(\varepsilon) \triangleq \begin{cases} \gamma_k = 0 & \text{if } k \notin S \\ y_e^e = x_e \gamma + \varepsilon^e & \text{where } \varepsilon^e \sim F_\varepsilon \forall e \in S \end{cases}$$

- * if $\exists \gamma \in R^P$ such that $H_{0,r,S}(\varepsilon)$ is true
then S : plausible causal predictors

* Define $S(\varepsilon) = \bigcap_{S, H_{0,r,S} \text{ is true}} S$

$S(\varepsilon) \triangleq$ Identifiable causal predictors

Easy to notice :

① Since $H_{0,s^*,r}(\varepsilon)$ is true.

$$\bigcap_{S, H_{0,r,S}} S \subseteq S^*$$

$$\therefore S(\varepsilon) \subseteq S^*$$

② If $\varepsilon_1 \subseteq \varepsilon_2$, $S(\varepsilon_1) \subseteq S(\varepsilon_2) \subseteq S^*$
 \implies closer to identifying S^*

* $\Gamma_S(\mathcal{E}) \triangleq \{\gamma \in \mathbb{R}^P \mid H_0, \text{r.s. is true}\}$

» $\Gamma(\mathcal{E}) \triangleq \bigcup_S \Gamma_S(\mathcal{E})$

Easy to notice :

if $\mathcal{E}_1 \subseteq \mathcal{E}_2$, $\Gamma(\mathcal{E}_2) \subseteq \Gamma(\mathcal{E}_1)$

Alternate definition of null hypothesis :

$\exists \beta \in \mathbb{R}^P$, such that $\forall e \in \mathcal{E}, \beta = \beta^{pred, e}(s)$

$$\beta^{pred, e}(s) = \underset{\beta, \beta_k=0, k \neq s}{\operatorname{argmin}} \left\{ E(Y_e - \beta x^e)^2 \right\}$$

Identifiability Results :

formalizing interventions : A_e : contains indices
of variables that we intervene on

$$A_1 = \emptyset$$

$$\beta_{jik}^e = \begin{cases} \beta_{jk}^i & \text{if } j \notin A_e \\ 0 & \text{if } j \in A_e \end{cases} \quad \rightarrow D_o \text{ interventions}$$

$$\varepsilon_j^e = \begin{cases} \varepsilon_j^i & \text{if } j \notin A_e \\ a_j^e & \text{if } j \in A_e \end{cases}$$

$$\beta_{jk}^e = \beta_{jk}^i + j$$

$$\varepsilon_j^e = \begin{cases} \varepsilon_j^i & \text{if } j \notin A_e \\ A_j^e \cdot \varepsilon_j^i & \text{if } j \in A_e \\ (0, \lambda) & C_j^e + \varepsilon_j^i \text{ if } j \in A_e \end{cases}$$

\rightarrow Noise intervention

Simultaneous: pool all intervention environments into $e = 2$
 + noise intervention as here

Identifiability Result:

Theorem 2 Consider a (linear) Gaussian SEM as in (19) and (20) with interventions. Then, with $S(\mathcal{E})$ as in (6), all causal predictors are identifiable, that is

$$S(\mathcal{E}) = \mathbf{PA}(Y) = \mathbf{PA}(1) \quad (22)$$

if one of the following three assumptions is satisfied:

- i) The interventions are **do-interventions** (Section 4.2.1) with $a_j^e \neq E(X_j^1)$ and there is at least one single intervention on each variable other than Y , that is for each $j \in \{2, \dots, p+1\}$ there is an experiment e with $\mathcal{A}^e = \{j\}$.
- ii) The interventions are **noise interventions** (Section 4.2.2) with $1 \neq E(A_j^e)^2 < \infty$, and again, there is at least one single intervention on each variable other than Y . If the interventions act additively rather than multiplicatively, we require $EC_j^e \neq 0$ or $0 < \text{Var } C_j^e < \infty$.
- iii) The interventions are **simultaneous noise interventions** (Section 4.2.3). This result still holds if we allow changing linear coefficients $\beta_{j,k}^{e=2} \neq \beta_{j,k}^{e=1}$ in (21) with (possibly random) coefficients $\beta_{j,k}^{e=2}$.

The statements remain correct if we replace the null hypothesis (10) with its weaker version (16).

These are examples for sufficient conditions for identifiability but there may be many more. For example, one may also consider random coefficients or changing graph structures (only the parents of Y must remain the same).

Extensions :

Some cases these strong conditions can be relaxed

Analysis of Estimation

Assume we have level- α test available

by defn: $\sup_P P(H_{0,S}(\varepsilon) \text{ rejected} | H_{0,S}(\varepsilon) \text{ is true}) \leq \alpha$

$\therefore S(\varepsilon) \subseteq S^*$ with $P \geq 1-\alpha$

$\gamma^* \in \hat{\Gamma}(\varepsilon)$ with $P \geq 1-2\alpha$

Proof:

$$P(S^* \text{ rejected}) \leq \alpha$$

$$\Rightarrow P(S^* \text{ accepted}) \geq 1-\alpha \Rightarrow P(S(\varepsilon) \subseteq S^*) \geq 1-\alpha$$

Note: $\hat{\Gamma}(\varepsilon) = \bigcup \hat{C}(S) : S \text{ makes null hypothesis true}$

$\hat{C}(S)$: $1-\alpha$ confidence set for θ or
 $\theta^{pred}(S)$

$$P(\gamma^* \notin \hat{\Gamma}(\varepsilon)) = P(S^* \text{ rejected}) + P(\gamma^* \notin \hat{C}(S))$$

$\uparrow \alpha$ $\hookrightarrow \alpha$

$$P(\gamma^* \notin \hat{\Gamma}(\varepsilon)) \leq 2\alpha$$

$$P(\gamma^* \in \hat{\Gamma}(\varepsilon)) \geq 1-2\alpha$$

Generic method for invariant prediction

- 1) For each set $S \subseteq \{1, \dots, p\}$, test whether $H_{0,S}(\mathcal{E})$ holds at level α (we will discuss later concrete examples).
- 2) Set $\hat{S}(\mathcal{E})$ as

$$\hat{S}(\mathcal{E}) := \bigcap_{S: H_{0,S}(\mathcal{E}) \text{ not rejected}} S. \quad (12)$$

- 3) For the confidence sets, define

$$\hat{\Gamma}(\mathcal{E}) := \bigcup_{S \subseteq \{1, \dots, p\}} \hat{\Gamma}_S(\mathcal{E}), \quad (13)$$

where

$$\hat{\Gamma}_S(\mathcal{E}) := \begin{cases} \emptyset & H_{0,S}(\mathcal{E}) \text{ can be rejected at level } \alpha \\ \hat{C}(S) & \text{otherwise.} \end{cases} \quad (14)$$

Here, $\hat{C}(S)$ is a $(1 - \alpha)$ -confidence set for the regression vector $\beta^{\text{pred}}(S)$ that is obtained by pooling the data.

Last piece : Coming up with a hypothesis
test at level alpha .

Method II: Invariant prediction using fast(er) approximate test on residuals

- 1) For each $S \subseteq \{1, \dots, p\}$ and $e \in \mathcal{E}$:
 - (a) Fit a linear regression model on all data to get an estimate $\hat{\beta}^{\text{pred}}(S)$ of the optimal coefficients using set S of variables for linear prediction in regression. Let $R = Y - X\hat{\beta}^{\text{pred}}(S)$.
 - (b) Test the null hypothesis that the mean of R is identical for each set I_e and $e \in \mathcal{E}$, using a two-sample t-test for residuals in I_e against residuals in I_{-e} and combining via Bonferroni correction across all $e \in \mathcal{E}$. Furthermore, test whether the variances of R are identical in I_e and I_{-e} , using an F-test, and combine again via Bonferroni correction for all $e \in \mathcal{E}$. Combine the two p -values of equal variance and equal mean by taking twice the smaller of the two values. If the p -value for the set S is smaller than α , we reject the set S .
- 2) As in the generic algorithm, using (12).
- 3) If we do not reject a set S we set $\hat{\Gamma}_S(\mathcal{E}) = \emptyset$. Otherwise, we set $\hat{\Gamma}_S(\mathcal{E})$ to be the conventional $(1 - \alpha)$ -confidence region for $\beta^{\text{pred}}(S)$ when using all data simultaneously. For simplicity, we will use rectangular confidence regions, exactly as in step 3 of Method I.

2 things not covered:

Extension to non-linear models
Extension to include hidden variables

} Basic extension allowed

Model misspecification: what happens if

we use linear model on non-linear SEM?

→ If S^* is accepted : then confidence statements hold
or satisfies hypothesis

→ If nothing satisfies : $\hat{S}(\varepsilon) = \emptyset$
max power to reject false hypothesis

→ If S^* does not, but something else : Under mild assump.

$\hat{S}(\varepsilon) = AN(Y)$
as opposed to $PA(Y)$

Invariant Risk Minimization

Invariant Representations of objects relate to the causal explanation of the object itself.

Canonical Example : Cows on grass

Goal: Enable OOD Generalization by learning invariant predictors

Approach: Learn a data representation such that it elicits an optimal predictor on top of it that is invariant

Notation Setup:

R : risk, R^e : risk when expectation is considered over $e \in E$

E_{all} : all environments in a modal realism sense

what's fair game? SEM(y): Can only be varied by varying noise model's variance in a finite range

SEM (anything else): mess arbitrarily

Note: ICP requires $P(Y|Pa(Y))$ invariant
IRM requires $E[y|Pa(y)]$ invariant

$$R^{OOD}(f) \triangleq \max_{e \in E_{\text{all}}} R^e(f) \quad \xrightarrow{\text{Goal:}} \quad f = \underset{f}{\operatorname{argmin}} \quad R^{OOD}(f)$$

$$R^{rob}(f) \triangleq \max_{e \in E_{tr}} R^e(f)$$

↓ Set of training environments
[we have access to labeled samples]

Some candidate approaches & their shortcomings

1) ERM : obvious shortcomings discussed & motivated ERM
(pooling data is ERM)

2) $\min R^{\text{rob}}(f)$: can be shown is weighted ERM under some regularity conditions

3) Domain Adaptation strategy : we are okay with $P(X_1), P(X_2)$,
 $P(Y)$, etc to change.

not justified to enforce $\mathbb{E}(X_1, X_2)$ with invariant distribution -

4) ICP : allow noise model to change variance in finite range

IRM :

To learn $f: \mathcal{X} \rightarrow \mathcal{Y}$, learn $\Phi: \mathcal{X} \rightarrow H$ and $w: H \rightarrow \mathcal{Y}$, $f = w \circ \Phi$.

key idea: w on top of Φ should be simultaneously optimal for all $e \in \mathcal{E}_{\text{all}}$

why? Common risks elicit optimal predictors of the form $E[y^e | x^e]$

$$\begin{aligned}\because \text{Invariance of } w \Rightarrow E[y^e | \Phi(x^e) = z] \\ = E[y^{e'} | \Phi(x^{e'}) = z]\end{aligned}$$

formally :

$$\Phi: \mathcal{X} \rightarrow H, w: H \rightarrow \mathcal{Y}$$

objective: $\min_{\Phi, w} \sum_{e \in \mathcal{E}_{\text{tr}}} R^e(w \circ \Phi)$. such $w \in \arg \min_{w: H \rightarrow \mathcal{Y}}$ that $R^e(w \circ \Phi) + e \in \mathcal{E}$

IRM:

IRM in its original form: Complex bi-level optimisation problem

Several relaxations:

a) Make constraint soft

$$\min_{e \in E^{\text{tr}}} R^e(w \circ \Xi) + \gamma D(w, \Xi, e)$$

choice of D :

Say $w \circ \Xi$ linear (w is LS Regression on $\Xi(x)$)

$$w_{\Xi}^e = \underset{x^e}{E} \left[\Phi(x^e) \Xi(x^e)^T \right]^{-1} \underset{x^e, y^e}{E} [\Xi(x^e) y^e]$$

$$D_{\text{dist}}(w, \Xi, e) = \|w - w_{\Xi}^e\|^2$$

poor: as $\Xi \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & c \end{bmatrix}$, as $c \rightarrow 0$, $\Xi \rightarrow \Xi^*$.

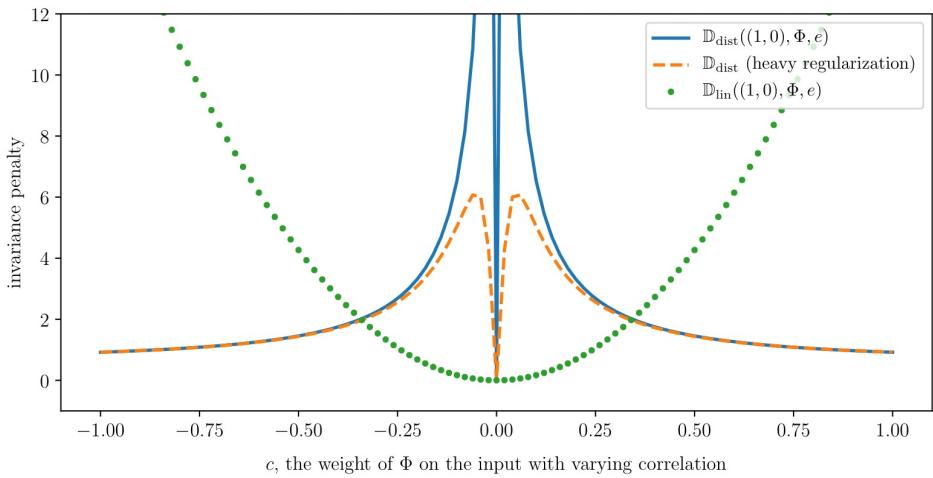
but as $c \rightarrow 0$, $D_{\text{dist}} \rightarrow \infty$

similarly as $c \rightarrow \infty$, $D_{\text{dist}} \rightarrow 0$

Root of this problem: Inverse

fix it: $D_{\text{Lin}}(w, \Xi, e)$

$$= \left\| \underset{x^e}{E} [\Phi(x^e) \Xi(x^e)^T] w - \underset{x^e, y^e}{E} [\Xi(x^e) y^e] \right\|$$



New Issue: Consider same mapping $w \circ \Phi$

$$\text{but } w \circ \tilde{\Phi} = \underbrace{w \circ \psi^{-1}}_{\tilde{w}} \circ \psi \circ \tilde{\Phi}$$

ERM part is the same

but say $\psi = \gamma$, $\gamma \in (0, 1)$

then $D'_{\text{Lin}} = \gamma D_{\text{Lin}} \rightarrow 0$ as $\gamma \rightarrow 0$

Final fix: fixing the linear classifier

Rephrase IRM: find a data representation $\tilde{\Phi}$ such that optimal linear classifier on top of that is \tilde{w} .

So optimization only over Φ

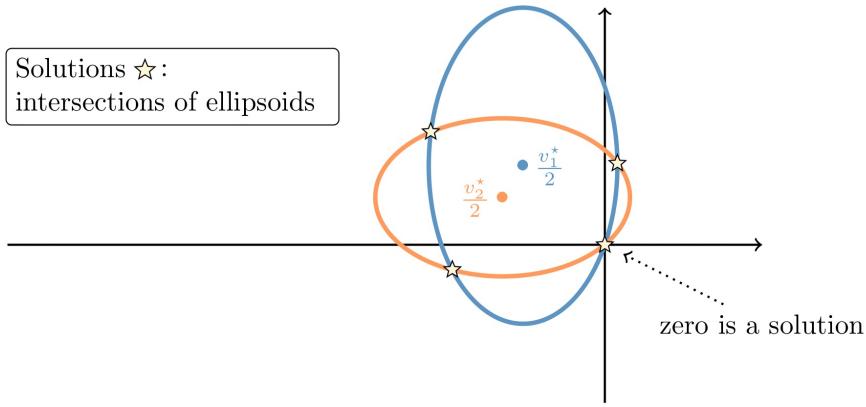
why?

Theorem 4. For all $e \in \mathcal{E}$, let $R^e : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex differentiable cost functions. A vector $v \in \mathbb{R}^d$ can be written $v = \Phi^\top w$, where $\Phi \in \mathbb{R}^{p \times d}$, and where $w \in \mathbb{R}^p$ simultaneously minimize $R^e(w \circ \Phi)$ for all $e \in \mathcal{E}$, if and only if $v^\top \nabla R^e(v) = 0$ for all $e \in \mathcal{E}$. Furthermore, the matrices Φ for which such a decomposition exists are the matrices whose nullspace $\text{Ker}(\Phi)$ is orthogonal to v and contains all the $\nabla R^e(v)$.

So, any linear invariant predictor can be decomposed as linear data representations of different ranks. In particular, we can restrict our search to matrices $\Phi \in \mathbb{R}^{1 \times d}$ and let $\tilde{w} \in \mathbb{R}^1$ be the fixed scalar 1.0. This translates (5) into:

$$L_{\text{IRM}, w=1.0}(\Phi^\top) = \sum_{e \in \mathcal{E}_{\text{tr}}} R^e(\Phi^\top) + \lambda \cdot \mathbb{D}_{\text{lin}}(1.0, \Phi^\top, e). \quad (6)$$

Solutions \star :
intersections of ellipsoids



Finally:

$$\text{IRM VI: } \min_{\Phi} \sum_{e \in \mathcal{E}_{\text{tr}}} R^e(\Phi) + \lambda \| \nabla_{w, w=1.0} R^e(w \Phi) \|^2$$

Generalization theory for IRM

when does invariance on Σ_{tr} translate to invariance on Σ_{all} ?

Linear General position of degree r :

$$|\Sigma_{\text{tr}}| > d-r + \frac{d}{r} \text{ for some } r \in \mathbb{N} \text{ and}$$

$$\forall x \neq 0 \in \mathbb{R}^d$$

$$\dim(\text{span}\left(\left\{\underset{x \in \Sigma}{E}[x x^T]x - \underset{x \in \Sigma}{E}[x^T e]\right\}\right)) > d-r$$

then if $y^e = z_1^e + e^e$, $z_1^e \perp e^e$, $E[e^e] = 0$

$x^e = s(z_1, z_2)$, s is invertible in part

$$(ii) \tilde{s}(s(z_1, z_2)) = z_1$$

then if $\Phi : \mathbb{R}^{d \times d}$, rank $r > 0$,

if $|\Sigma_{\text{tr}}| > d-r + \frac{d}{r}$ and lie in L.G.P of

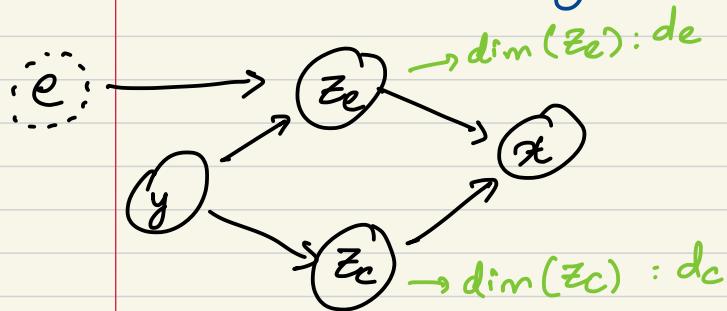
degree r , $\Phi E_{x^e} [x^e x^e]^T \Phi^T$

$$- \Phi E_{x^e y^e} [x^e y^e] = 0 \text{ on } \Sigma_{\text{tr}}$$

if and only if $\Phi^T w$ is invariant on Σ_{all}

The Risks of Invariant Risk Minimization

Informal results based on the following pedagogical model (which can be relaxed in ways)



$$R^e(\Phi, \hat{\beta}) = E_{(x,y) \sim p^e} [l(g(\hat{\beta}^\top \phi(x)), y)]$$

(I) linear f

- $|\varepsilon| > de$: IRM is dope
- $|\varepsilon| \leq de$: IRM solutions $\Phi, \hat{\beta}$ relies upon non-invariant features

(II) linear f : ∃ linear predictor $\Phi, \hat{\beta}$ which uses only environmental features yet achieves lower risk than optimal invariant predictor

(III) Non-linear f : on train: nearly identical to OIP

on test : equivalent to IRM
if test correlations are reversed
compared to train : this achieves
almost 0 accuracy

Very rough pic

$$\text{IRM} \stackrel{\text{train}}{=} \text{Middle guy} \stackrel{\text{test}}{=} \text{ERM}$$

if test highly opposite to
train

IRM gets destroyed !

TABLE-1: A snapshot of results from Gulrajani et al, showing AVERAGE OOD accuracies for ERM and IRM across the seven datasets Considered with three model Selection criteria.

Model selection method: training domain validation set								
Algorithm	CMNIST	RMNIST	VLCS	PACS	Office-Home	TerraInc	DomainNet	Avg
ERM	52.0 ± 0.1	98.0 ± 0.0	77.4 ± 0.3	85.7 ± 0.5	67.5 ± 0.5	47.2 ± 0.4	41.2 ± 0.2	67.0
IRM	51.8 ± 0.1	97.9 ± 0.0	78.1 ± 0.0	84.4 ± 1.1	66.6 ± 1.0	47.9 ± 0.7	35.7 ± 1.9	66.0

Model selection method: Leave-one-domain-out cross-validation								
Algorithm	CMNIST	RMNIST	VLCS	PACS	Office-Home	TerraInc	DomainNet	Avg
ERM	34.2 ± 1.2	98.0 ± 0.0	76.8 ± 1.0	83.3 ± 0.6	67.3 ± 0.3	46.2 ± 0.2	40.8 ± 0.2	63.8
IRM	36.3 ± 0.4	97.7 ± 0.1	77.2 ± 0.3	82.9 ± 0.6	66.7 ± 0.7	44.0 ± 0.7	35.3 ± 1.5	62.9

Model selection method: Test-domain validation set (oracle)								
Algorithm	CMNIST	RMNIST	VLCS	PACS	Office-Home	TerraInc	DomainNet	Avg
ERM	58.5 ± 0.3	98.1 ± 0.1	77.8 ± 0.3	87.1 ± 0.3	67.1 ± 0.5	52.7 ± 0.2	41.6 ± 0.1	68.9
IRM	70.2 ± 0.2	97.9 ± 0.0	77.1 ± 0.2	84.6 ± 0.5	67.2 ± 0.8	50.9 ± 0.4	36.0 ± 1.6	69.2

TABLE-2: Domain generalization accuracies on CMNIST per domain.

Model selection method: training domain validation set								
Algorithm	0.1	0.2	0.9					
ERM	72.7 ± 0.2	73.2 ± 0.3	10.0 ± 0.0					
IRM	72.0 ± 0.3	73.2 ± 0.0	10.1 ± 0.2					

Model selection method: leave-one-domain-out cross-validation								
Algorithm	0.1	0.2	0.9					
ERM	46.0 ± 3.4	46.6 ± 3.8	10.0 ± 0.1					
IRM	49.3 ± 0.9	49.5 ± 0.2	10.0 ± 0.2					

Model selection method: test-domain validation set (oracle)								
Algorithm	0.1	0.2	0.9					
ERM	72.3 ± 0.6	73.1 ± 0.3	30.0 ± 0.3					
IRM	72.7 ± 0.1	72.8 ± 0.3	65.2 ± 0.8					

→ The CMNIST Dataset

- Assign a preliminary binary label \hat{y} to the image based on the digit:
 $\hat{y} = 0$ for digits 0-4 and $\hat{y} = 1$ for 5-9.

- Obtain the final label \hat{y} by flipping y with probability 0.25
- Sample the colour ID z by flipping y with probability p^e , where $p^e = 0.2$ in the first environment, 0.1 in the second, and 0.9 in the third.
- Colour the image red if $z=1$ or green if $z=0$.

IRM vs ERM: A SAMPLE COMPLEXITY PERSPECTIVE.

→ Setup:

- Dataset $D = \{D^e\}_{e \in E_{tr}}$, which is a collection of datasets $D^e = \{(x_i^e, y_i^e, e)\}_{i=1}^{n_e}$ obtained from a set of training environments E_{tr} .
- Probability distribution $\{\pi^e\}_{e \in E_{tr}}$, where π^e is the probability that a training data point is from environment e .
- Define a predictor $f: X \rightarrow \mathbb{R}$ and the space \mathcal{F} of all possible maps from $X \rightarrow \mathbb{R}$.
- Define risk achieved by f in environment e as $R^e(f) = \mathbb{E}^e [l(f(x^e), y^e)]$, where l is the loss. Overall expected risk across the training environments is $R(f) = \sum_{e \in E_{tr}} \pi^e R^e(f)$.

→ OOD Generalization Problem:

$$\min_{f \in \mathcal{F}} \max_{e \in E_{all}} R^e(f). \rightarrow \textcircled{1}$$

→ Assumption -1 : Invariance Condition.

There exists a representation $\underline{\Phi}^*$ that transforms X^e to $Z^e = \underline{\Phi}^*(x^e)$ and $\forall e, o \in E_{\text{all}}, \forall z \in \underline{\Phi}^*(x)$

satisfies $E^e[Y^e | Z^e = z] = E^o[Y^o | Z^o = z]$. Also, $\forall e \in E_{\text{all}}$, $\forall z \in \underline{\Phi}^*(x)$, $\text{Var}^e[Y^e | Z^e = z] = \xi^2$, where Var^e is the conditional variance.

→ Using assumption-1, define an invariant map $m : \underline{\Phi}^*(X) \rightarrow \mathbb{R}$

as :

$$\forall z \in \underline{\Phi}^*(X), m(z) = E^e[Y^e | Z^e = z], \quad \rightarrow \textcircled{2}$$

(Since $Z^e = \underline{\Phi}^*(x^e)$).

→ Assumption -2 : There exists an environment $e \in E_{\text{all}}$ s.t.

$$Y^e \perp X^e | Z^e$$

i.e., in environment e , the information that X^e has about Y^e is also contained in Z^e .

→ Define composition $m \circ \underline{\Phi}^*$, $\forall x \in X, m \circ \underline{\Phi}^*(x) = E^e[Y^e | Z^e = \underline{\Phi}^*(x)]$

→ Proposition-1 : If l is the square loss, and assumptions 1 and 2 hold, then $m \circ \underline{\Phi}^*$ solves the OOD problem.

→ Recall the IRM objective :

- A data representation $\underline{\Phi}$ elicits an invariant predictor $w \circ \underline{\Phi}$ across environments $e \in E_{\text{tar}}$ if there is a classifier w that achieves minimum risk simultaneously for all environments.

- IRM selects the invariant predictors with the least sum of risk across all training environments:

$$\min_{\substack{\Phi \in \mathcal{H}_\Phi \\ \omega \in \mathcal{H}_\omega}} R(\omega \cdot \Phi) = \sum_{e \in E_{tr}} \pi^e R^e(\omega \cdot \Phi) \quad \rightarrow \textcircled{3}$$

$$\text{s.t. } \omega \in \operatorname{argmin}_{\bar{\omega} \in \mathcal{H}_\omega} R^e(\bar{\omega} \cdot \Phi), \forall e \in E_{tr}.$$

→ $\omega \cdot \Phi^*$ is a feasible solution $\textcircled{3}$ and is also the ideal Solution as it solves $\textcircled{1}$.

→ If gradient constrained formulation of IRM:

$$\min_{\Phi \in \mathcal{H}_\Phi} R(\Phi) \quad \rightarrow \textcircled{4}$$

$$\text{s.t. } \|\nabla_{\omega} |_{\omega=1.0} h^e(\omega \cdot \Phi)\| = 0, \forall e \in E_{tr}.$$

→ ϵ -approximation of the IRM objective and its empirical version.

- Define $R'(\Phi) = \sum_{e \in E_{tr}} \pi^e \|\nabla_{\omega} |_{\omega=1.0} R^e(\omega \cdot \Phi)\|^2$ and set

$$S^v(\epsilon) = \{\Phi \mid R'(\Phi) \leq \epsilon, \Phi \in \mathcal{H}_\Phi\}$$

- The ϵ -approximation of $\textcircled{4}$ is:

$$\min_{\Phi \in S^v(\epsilon)} R(\Phi). \rightarrow \textcircled{5}$$

- An empirical version of $\textcircled{5}$ replaces R and R' with empirical estimators \hat{R} and \hat{R}' respectively:

$$\min_{\Phi \in \hat{S}^v(\epsilon)} \hat{R}(\Phi) \rightarrow \textcircled{6}$$

$$\text{where } \hat{S}^v(\epsilon) = \{\Phi \mid \hat{R}'(\Phi) \leq \epsilon, \Phi \in \mathcal{H}_\Phi\}$$

\hat{R} comes from the sample mean of loss across samples from D .

\hat{R}' is obtained in Supplement 7.3.1 (pages 23 and 24).

→ Proposition-2 (highlight):

Sample complexity grows as $\mathcal{O}\left(\max\left\{\frac{1}{\epsilon^2}, \frac{1}{\delta^2}\right\}\right)$

to ensure that every solution to empirical IRM (⑥)
is a γ -approximation of IRM with probability at
least $1-\delta$

→ Proposition-3 (Shalev-Shwartz and Ben-David, 2014):

Sample complexity grows as $\mathcal{O}\left(\frac{1}{\gamma^2}\right)$ to ensure
that every solution to ERM is a γ -approximation
of expected risk minimization with probability at
least $1-\delta$.

→ OOD PERFORMANCE (ERM vs IRM) : COVARIATE SHIFT.

→ Assumption-4: Invariance w.r.t all features.

$\forall e, o \in E_{all}$ and $\forall x \in \mathcal{X}$, $E[Y^e | X^e = x] = E[Y^o | X^o = x]$

$\forall e \in E_{all}$, $X^e \sim P_{X^e}^e$ and the Support of $P_{X^e}^e$ is equal to \mathcal{X} .

→ The first part of the assumption is equivalent to setting $\Phi^* = I$
in assumption-1. If $\Phi^* = I$, then m in equation ②

Simplifies to $m(u) = E[Y^e | X^e = u]$ and solves the OOD
problem ①

→ Proposition-4: Let L be the square loss, and

$\tilde{\lambda} = \min_{\Phi_1, \Phi_2 \in \mathcal{H}_{\Phi}} |R(\Phi_1) - R(\Phi_2)|$ i.e., the minimum
separation between the risks of any two distinct

hypotheses in \mathcal{H}_{Φ} .

For every $\gamma > 0$, $\epsilon > 0$, and $\delta \in (0, 1)$ if \mathcal{H}_{Φ} is
a finite hypothesis class, $m \in \mathcal{H}_{\Phi}$, assumptions 3
and 4 hold, then:

- If the number of samples $|D|$ is greater than $\max \left\{ \frac{8L^2}{\gamma^2} \log \left(\frac{4|H_{\Phi}|}{\delta} \right), \frac{16L^4}{\epsilon^2} \log \frac{2}{\delta} \right\}$, then with probability at least $1 - \delta$, every solution $\hat{\Phi}$ to EIRM (equation ⑥) satisfies $R(m) \leq R(\hat{\Phi}) \leq R(m) + 2\gamma$. If $2\gamma < \tilde{\lambda}$, then $\hat{\Phi} = m$.
- If the number of samples $|D|$ is greater than $\frac{8L^2}{\gamma^2} \log \left(\frac{2|H_{\Phi}|}{\delta} \right)$ then with probability at least $1 - \delta$, every solution $\hat{\Phi}^+$ to ERM satisfies $R(m) \leq R(\hat{\Phi}^+) \leq R(m) + 2\gamma$. If $2\gamma < \tilde{\lambda}$, then $\hat{\Phi}^+ = m$.

→ OOD GENERALIZATION (ERM vs IRM):
CONFFOUNDERS OR/AND ANTI-CAUSAL VARIABLES.

→ Recall the linear models from Arjovsky et al. (2019):

$$\begin{aligned} e &\sim \text{Categorical}(\{\pi^e\}_{0 \in E_{\text{tr}}}, \forall 0 \in E_{\text{tr}}, \pi^e > 0) \\ y^e &\leftarrow r^T(z_1^e) + \varepsilon^e, \quad \varepsilon^e \perp z_1^e, \quad \mathbb{E}[\varepsilon^e] = 0, \quad \mathbb{E}[(\varepsilon^e)^2] = \sigma^2 \\ x^e &\leftarrow s(z_1^e, z_2^e). \end{aligned}$$

→ Assumption-5

Z_1 component of S is invertible, i.e. $\exists \tilde{S}$ s.t.

$$\tilde{S}(s(z_1, z_2)) = Z_1, \text{ and } r \neq 0.$$

Define $\Sigma^e = \mathbb{E}[x^e x^{e T}]$. s.t. $\forall e \in E_{\text{tr}}$, Σ^e is positive definite.

The support of Z^e and norm of S are bounded.

→ If Z_2^e is an effect of y^e , then Z_2^e is an anti-causal variable.

→ If H^e causes both ε^e and Z_2^e , then H^e is a confounder.

→ Given the above, the linear model $\tilde{S}^T r$ solves the OOD problem in equation ① ($\tilde{\Phi}^* = \tilde{S}$, $m = r^T$ in proposition-1)

→ Assumption-6: Linear general position (from IRM)

→ Proposition-s (highlight):

If the slack on the IRM penalty (ϵ) is sufficiently small, and the data grows as $O(\frac{1}{\epsilon^2})$
 then every solution $\hat{\Phi}$ to EIRM is in $\sqrt{\epsilon}$ radius of the OOD solution, i.e.,
 $\|\hat{\Phi} - \tilde{\xi}^\top r\| = O(\sqrt{\epsilon})$

Table 1: Summary of (empirical) IRM vs. ERM for finite hypothesis class \mathcal{H}_Φ . ϵ : slack in IRM constraints, ν : approximation w.r.t optimal risk, δ : failure probability, \mathcal{E}_{tr} : set of training environments, n : data dimension, p : degree of the generative polynomial, L, L' : bound on loss & its gradients.

Assumptions	Method	Sample complexity	OOD
<i>Covariate shift case:</i> $\mathbb{E}^e[Y^e X^e]$ is invariant (Proposition 4)	ERM	$\frac{8L^2}{\nu^2} \log\left(\frac{2 \mathcal{H}_\Phi }{\delta}\right)$	Yes
	IRM	$\max\left\{\frac{8L^2}{\nu^2} \log\left(\frac{4 \mathcal{H}_\Phi }{\delta}\right), \frac{16L'^4}{\epsilon^2} \log\left(\frac{2}{\delta}\right)\right\}$	Yes
<i>Confounder/Anti-causal variable case:</i> $\mathbb{E}^e[Y^e \Phi^*(X^e)]$ is invariant, Linear, Polynomial models, $ \mathcal{E}_{tr} = \mathcal{O}(n^p)$ (Proposition 5, 6, 17)	ERM	$\frac{8L^2}{\nu^2} \log\left(\frac{2 \mathcal{H}_\Phi }{\delta}\right)$	No
	IRM	$\frac{16L'^4}{\epsilon^2} \log\left(\frac{2 \mathcal{H}_\Phi }{\delta}\right)$	Yes