
Domain Generalization :

Domain-invariant Representation Learning

Data Mining & Quality Analytics Lab.

2024. 01. 19

발표자: 정진용



발표자 소개



❖ 정진용 (Jinyong Jeong)

- 고려대학교 산업경영공학과 석·박사 통합과정(2021.09~)
- Data Mining & Quality Analytics Lab. (김성범 교수님)

❖ 관심 연구 분야

- Domain Generalization
- Semi-Supervised Learning & Class-Imbalanced Semi-Supervised Learning

❖ E-mail

- jy_jeong@korea.ac.kr

목차

1. Introduction

- Empirical Risk Minimization
- Background of Domain Generalization

2. Domain-invariant representation learning

- Invariant Risk Minimization

3. Conclusion



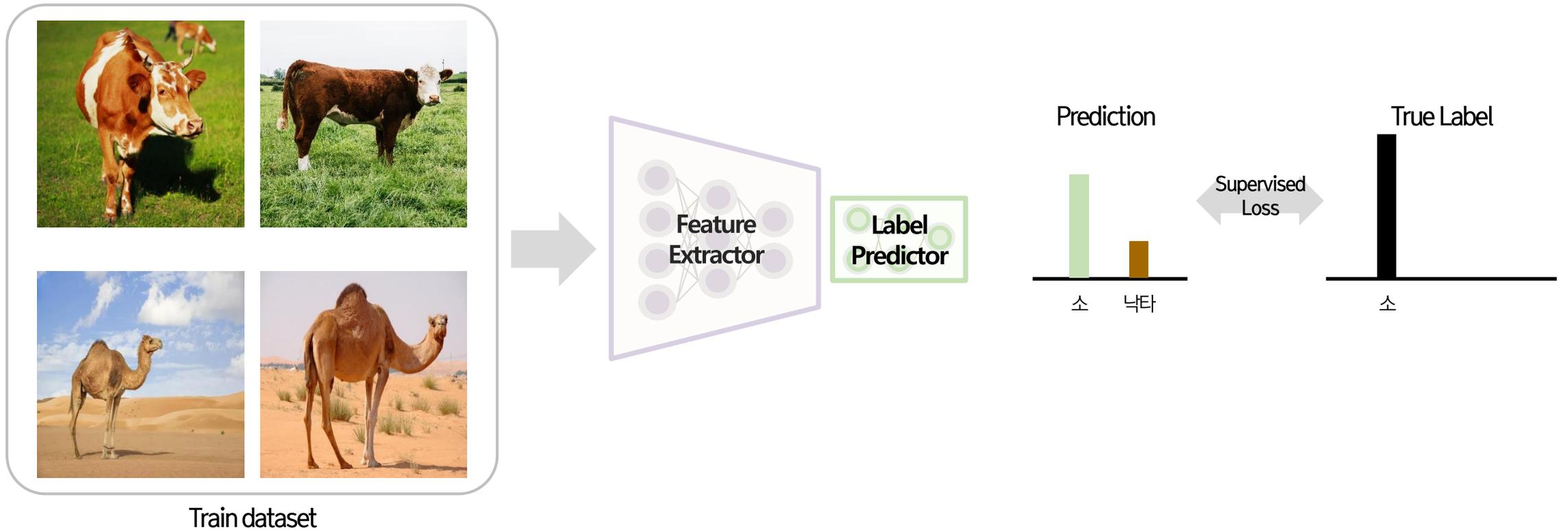
1. Introduction: Domain Generalization



Introduction

Empirical Risk Minimization

❖ 지도학습을 사용하여 분류 모델($f: X \rightarrow Y$)을 만들어보자



Introduction

Empirical Risk Minimization

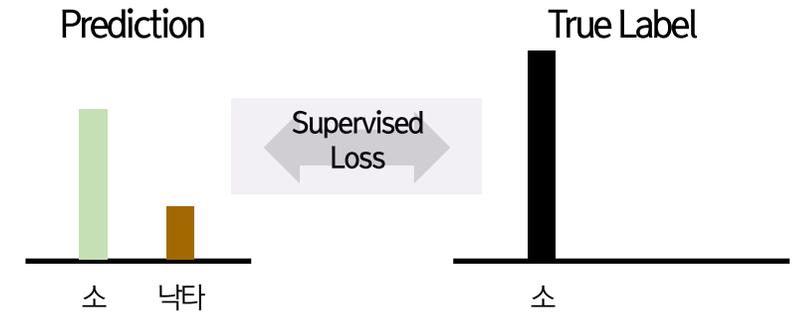
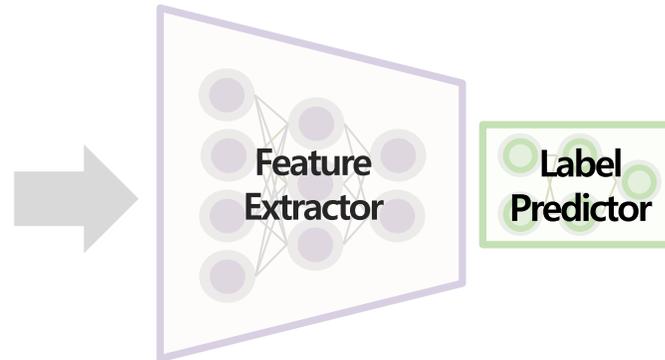
❖ 지도학습을 사용하여 분류 모델($f: X \rightarrow Y$)을 만들어보자

$$R(f) := \mathbb{E}_{(X,Y) \sim P(X,Y)} [L(f(X), Y)]$$

학습 목표 $f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)$



Train dataset



Introduction

Empirical Risk Minimization

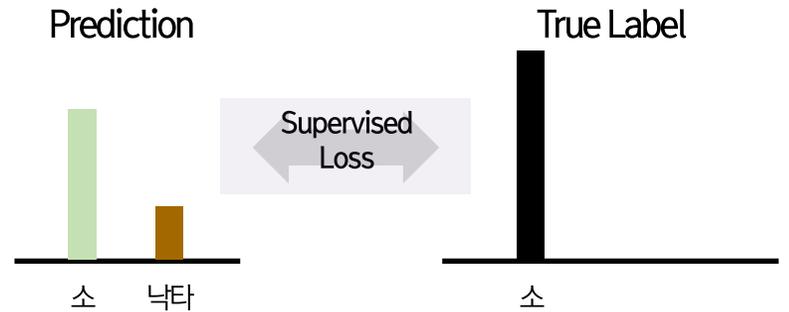
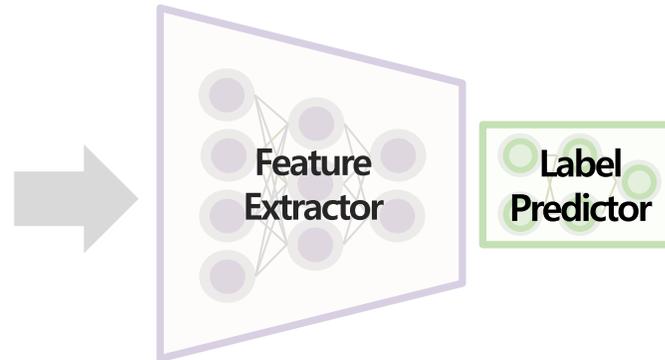
❖ 지도학습을 사용하여 분류 모델($f: X \rightarrow Y$)을 만들어보자

$$R(f) := \mathbb{E}_{(X,Y) \sim P(X,Y)} [L(f(X), Y)]$$

학습 목표 $f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)$



Train dataset



실제로는 X, Y 에 대한 확률 모분포 $P(X, Y)$ 를 알 수 없음



Introduction

Empirical Risk Minimization

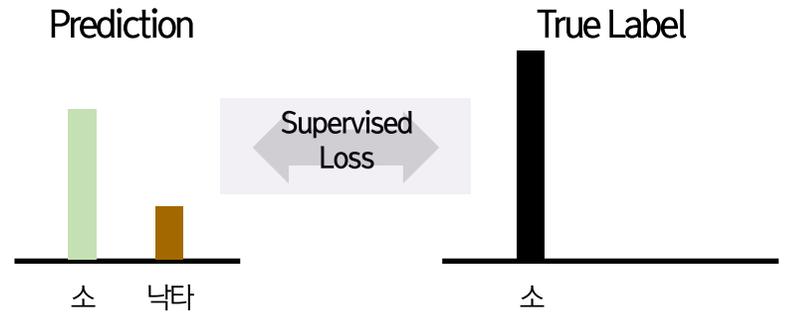
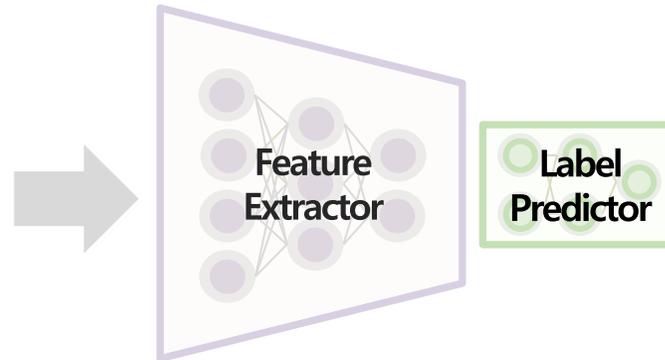
❖ 지도학습을 사용하여 분류 모델($f: X \rightarrow Y$)을 만들어보자

$$R(f) := \mathbb{E}_{(X,Y) \sim P(X,Y)} [L(f(X), Y)]$$

학습 목표 $f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)$



Train dataset



실제로는 X, Y 에 대한 확률 모분포 $P(X, Y)$ 를 알 수 없음
→ $P(X, Y)$ 에 속하는 데이터들을 샘플링

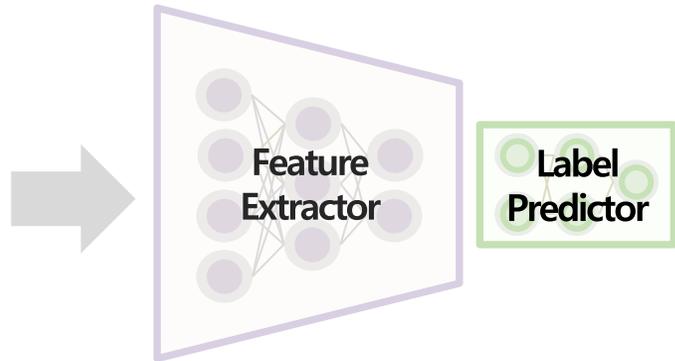
Introduction

Empirical Risk Minimization

❖ 지도학습을 사용하여 분류 모델($f: X \rightarrow Y$)을 만들어보자

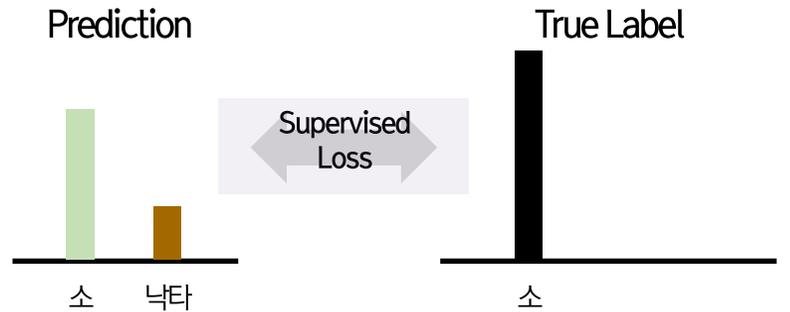


Train dataset



$$R(f) := \mathbb{E}_{(X,Y) \sim P(X,Y)} [L(f(X), Y)]$$

학습 목표 $f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)$



실제로는 X, Y 에 대한 확률 모분포 $P(X, Y)$ 를 알 수 없음
→ $P(X, Y)$ 에 속하는 데이터들을 샘플링

$$\text{Empirical } R(f) = \mathbb{E}_{X,Y \sim \hat{P}_{train}(X,Y)} [L(f(X), Y)] = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$



Introduction

Empirical Risk Minimization

❖ 학습된 분류 모델($f: X \rightarrow Y$)을 사용하여 추론을 해보자

$\hat{P}_{train}(X, Y)$ 에 의해서 학습된 f

$\hat{P}_{train}(X, Y) \approx \hat{P}_{test}(X, Y)$



Train dataset



Test dataset

Introduction

Empirical Risk Minimization

❖ 학습된 분류 모델($f: X \rightarrow Y$)을 사용하여 추론을 해보자

$\hat{P}_{train}(X, Y)$ 에 의해서 학습된 f



Train dataset

$\hat{P}_{test}(X, Y)$



Test dataset

Introduction

Empirical Risk Minimization

❖ 학습된 분류 모델($f: X \rightarrow Y$)을 사용하여 추론을 해보자

$\hat{P}_{train}(X, Y)$ 에 의해서 학습된 f



Train dataset

$\hat{P}_{test}(X, Y)$



Test dataset



Introduction

Empirical Risk Minimization

❖ 학습된 분류 모델($f: X \rightarrow Y$)을 사용하여 추론을 해보자

$\hat{P}_{train}(X, Y)$ 에 의해서 학습된 f



Train dataset



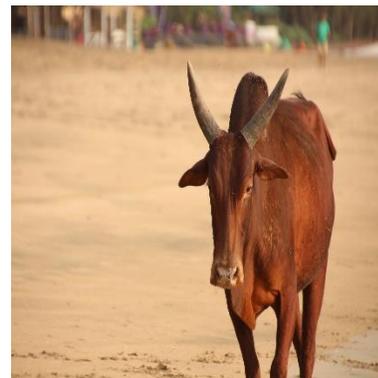
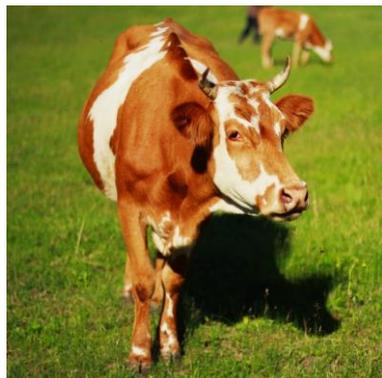
Test dataset

Introduction

Empirical Risk Minimization

❖ 학습된 분류 모델($f: X \rightarrow Y$)을 사용하여 추론을 해보자

$\hat{P}_{train}(X, Y)$ 에 의해서 학습된 f



Train dataset

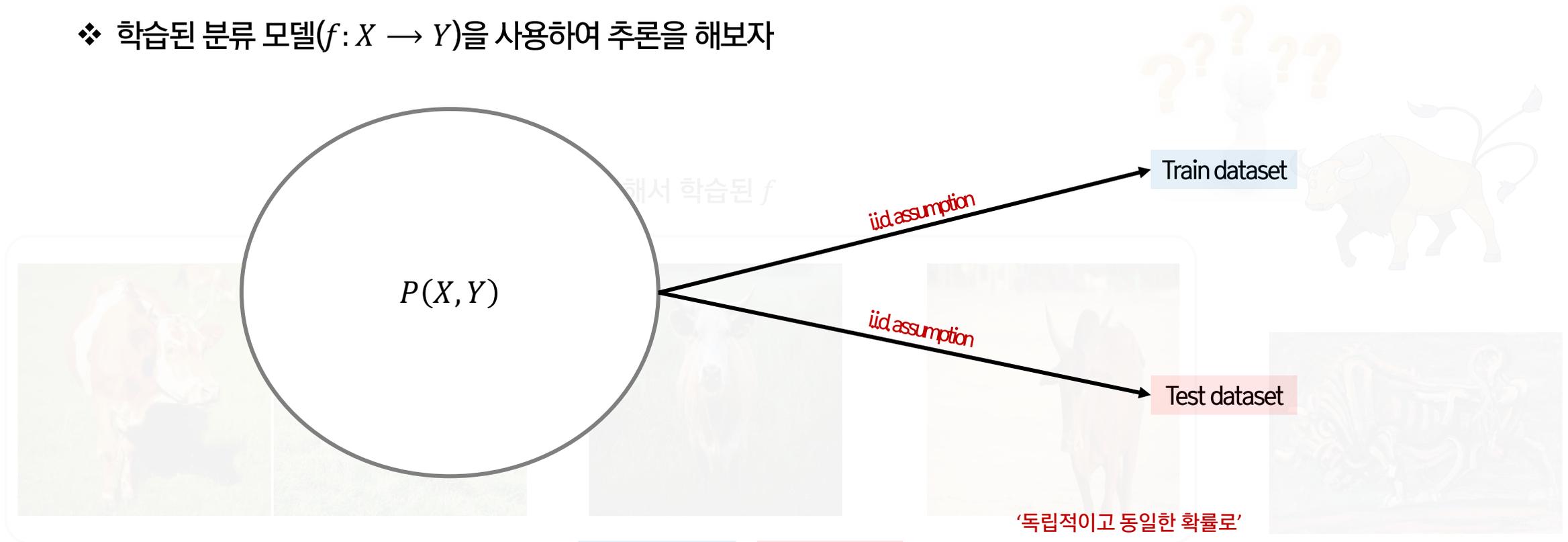


Test dataset

Introduction

Empirical Risk Minimization

❖ 학습된 분류 모델($f: X \rightarrow Y$)을 사용하여 추론을 해보자



기존 머신러닝 및 딥러닝 학습에서는 train dataset과 test dataset이 같은 분포에서 샘플링 되었다는 가정이 있음

→ 만약 iid 가정이 깨진다면, Empirical Risk Minimization을 통해 학습된 모델은 예측 성능이 저하될 수 있음

'Distribution shift'

Introduction

Background of Domain Generalization

❖ Distribution shift 상황에서도 모델 일반화 성능을 향상 시킬 수 있는 다양한 연구들이 존재

- Test data에 대한 성능을 높이자!

종료

Unsupervised domain adaptation Semi-supervised domain adaptation

도메인 적응
Domain adaptation

Source domain Target domain

2023년 09월 28일
Domain invariant model

Introduction to unsupervised domain adaptation

발표자: 이민정

📅 2022년 9월 16일
🕒 오후 1시 ~
📺 온라인 비디오 시청 (YouTube)

세미나 정보 보기 →

종료

How to Transfer Knowledge Across Domains by Deep Neural Network?

2022, 10, 28
Data Mining & Quality Analytics Lab.

How to Transfer Knowledge Across Domains

발표자: 김지현

📅 2022년 10월 28일
🕒 오후 1시 ~
📺 온라인 비디오 시청 (YouTube)

세미나 정보 보기 →

종료

Domain Generalization : How to improve the generalization ability of deep learning models?

DMQA Open Seminar (2023.07.21)
Data Mining & Quality Analytics Lab.

Domain Generalization: How to improve the generalization ability of deep learning models?

발표자: 김지현

📅 2023년 7월 21일
🕒 오후 12시 ~
📺 온라인 비디오 시청 (YouTube)

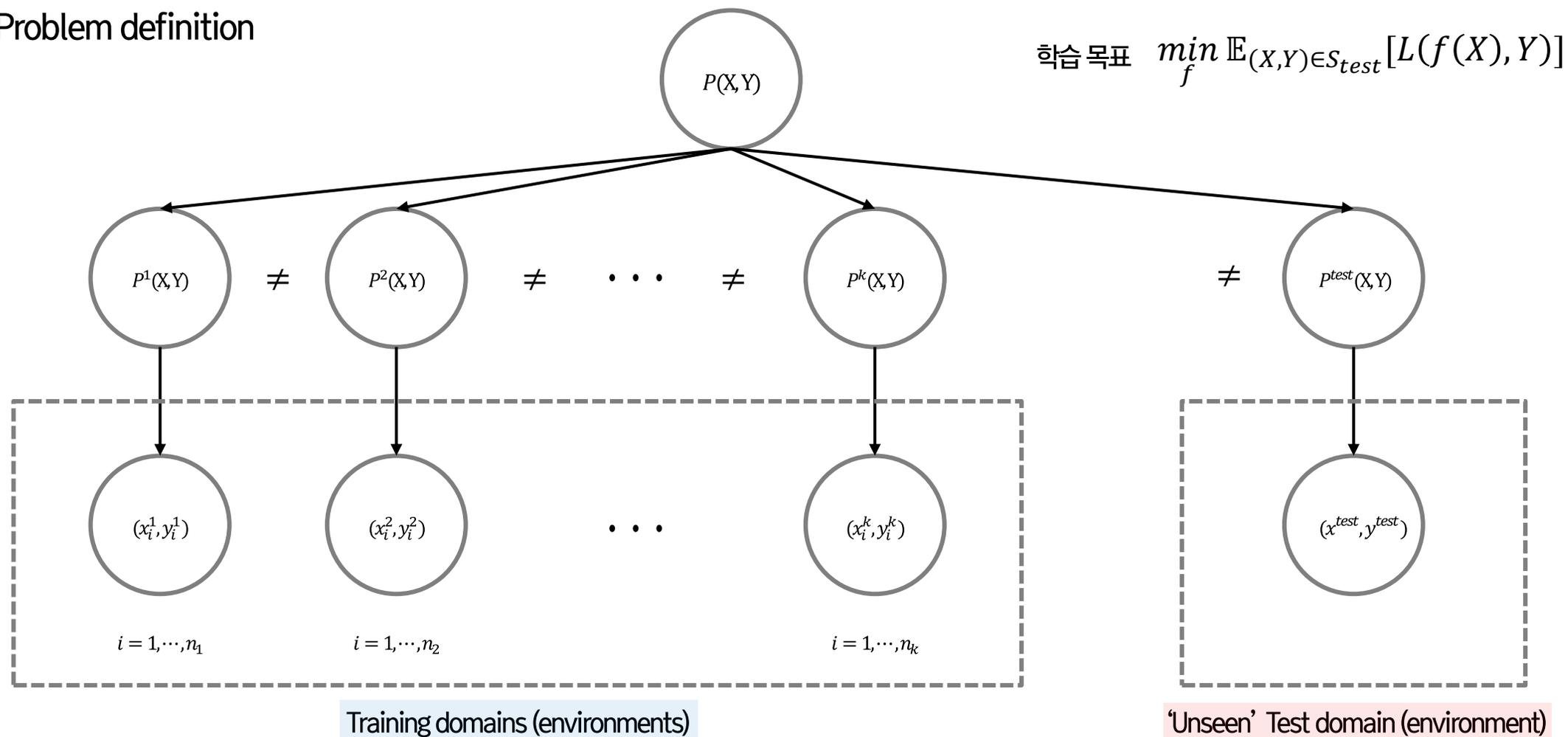
세미나 정보 보기 →



Introduction

Background of Domain Generalization

❖ Problem definition



Introduction

Background of Domain Generalization

❖ Taxonomy of domain generalization

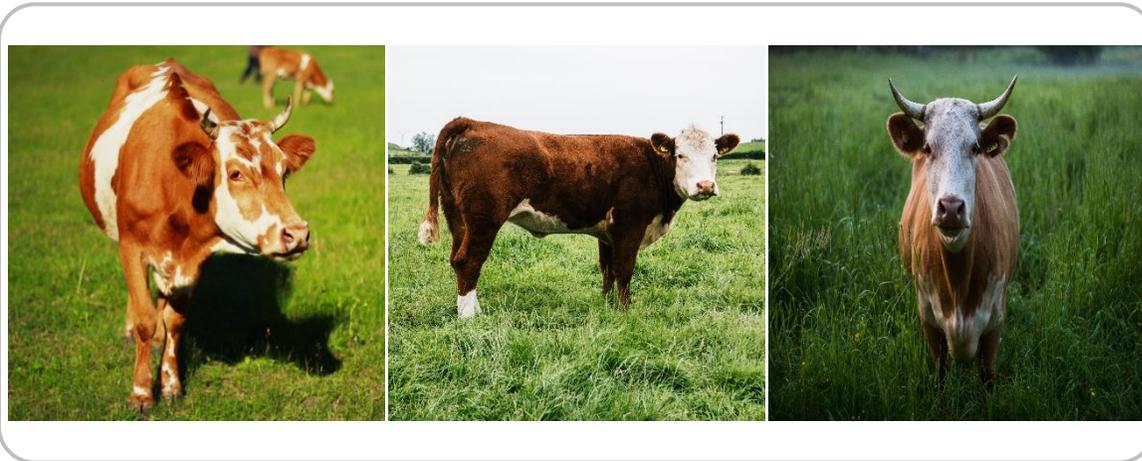


Domain-invariant representation learning

Invariant Risk Minimization

❖ Why is it a cow?

- ‘소’ 라는 것과 ‘green or glassy background’ 는 실제로 연관이 없음
 - Spurious correlation, not stable property
- ‘소의 형태’ 와 같은 invariant correlations을 배운다면 새로운 도메인에도 적용될 것임



Train domains



Test domains

Domain-invariant representation learning

Invariant Risk Minimization

❖ Invariant risk minimization (arxiv, 2019)

- Google Deepmind, Facebook AI, Stanford, Meta AI 연구원들에 의해 연구되었으며, 2024년 1월 19일 기준 1,681회 인용됨
- Domain-invariant representation learning에서 새로운 학습 패러다임을 제시한 논문

arXiv:1907.02893v3 [stat.ML] 27 Mar 2020

Invariant Risk Minimization

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, David Lopez-Paz

1 Introduction

Machine learning suffers from a fundamental problem. While machines are able to learn complex prediction rules by minimizing their training error, data are often marred by selection biases, confounding factors, and other peculiarities [49, 48, 23]. As such, machines justifiably inherit these data biases. This limitation plays an essential role in the situations where machine learning fails to fulfill the promises of artificial intelligence. More specifically, minimizing training error leads machines into recklessly absorbing all the correlations found in training data. Understanding which patterns are useful has been previously studied as a correlation-versus-causation dilemma, since spurious correlations stemming from data biases are unrelated to the causal explanation of interest [31, 27, 35, 52]. Following this line, we leverage tools from causation to develop the mathematics of spurious and invariant correlations, in order to alleviate the excessive reliance of machine learning systems on data biases, allowing them to generalize to new test distributions.

As a thought experiment, consider the problem of classifying images of cows and camels [4]. To address this task, we label images of both types of animals. Due to a selection bias, most pictures of cows are taken in green pastures, while most pictures of camels happen to be in deserts. After training a convolutional neural network on this dataset, we observe that the model fails to classify easy examples of images of cows when they are taken on sandy beaches. Bewildered, we later realize that our neural network successfully minimized its training error using a simple cheat: classify green landscapes as cows, and beige landscapes as camels.

To solve the problem described above, we need to identify which properties of the training data describe spurious correlations (landscapes and contexts), and which properties represent the phenomenon of interest (animal shapes). Intuitively, a correlation is spurious when we do not expect it to hold in the future in the same manner as it held in the past. In other words, spurious correlations do not appear to be stable properties [54]. Unfortunately, most datasets are not provided in a form amenable to discover stable properties. Because most machine learning algorithms depend on the assumption that training and testing data are sampled independently from the same distribution [51], it is common practice to shuffle at random the training and testing examples. For instance, whereas the original NIST handwritten data was collected from different writers under different conditions [19], the popular MNIST training and testing sets [8] were carefully shuffled to represent similar mixes of writers. Shuffling brings the training and testing distributions closer together, but



Domain-invariant representation learning

Invariant Risk Minimization

❖ Invariant predictor와 invariant correlation을 간단한 수식으로 확인

- Example : Structural equation model

$$\begin{array}{l}
 X_1 \stackrel{\text{noise}}{\leftarrow} \text{Gaussian}(0, \sigma^2) \\
 Y \leftarrow X_1 + \text{Gaussian}(0, \sigma^2) \\
 X_2 \leftarrow Y + \text{Gaussian}(0, 1) \quad \text{fixed noise}
 \end{array}$$

Noisy observation of Y ✓

$$e \in \mathcal{E}_{tr} = \{ \text{replace } \sigma^2 \text{ by } 10, \text{ replace } \sigma^2 \text{ by } 20 \}$$

$e=1$ $e=2$

이때, ~~environment~~ domain e에 대해서 Least squares predictor $\hat{Y}^e = X_1^e \hat{\alpha}_1 + X_2^e \hat{\alpha}_2$ 를 사용하여 Y를 예측하면?

- Regress from X_1^e : $\hat{\alpha}_1 = 1, \hat{\alpha}_2 = 0$
- Regress from X_2^e : $\hat{\alpha}_1 = 0, \hat{\alpha}_2 = \frac{\sigma(e)^2}{\sigma(e)^2 + 1}$
- Regress from (X_1^e, X_2^e) : $\hat{\alpha}_1 = 1/(\sigma(e)^2 + 1), \hat{\alpha}_2 = \sigma(e)^2/(\sigma(e)^2 + 1)$

$$Y^e = X_1^e + 0 \cdot X_2^e$$

regress from X_2

$$SSR = \sum (Y - \hat{d}_1 X_1 - \hat{d}_2 X_2)^2$$

$$\frac{\partial}{\partial \hat{d}_2} = -2 \sum (Y - \hat{d}_2 X_2) X_2 = 0$$

$$\hat{d}_2 = \frac{\sum X_2 Y}{\sum (X_2)^2} = \frac{\text{Cov}(X_2, Y)}{\text{Var}(X_2)}$$

$$\begin{aligned}
 \text{Cov}(X_2, Y) &= E[(X_2 - E[X_2])(Y - E[Y])] \\
 &= E[(X_2 - 0)(Y - 0)] \\
 &= E[Y^2 + \underbrace{(Y \times N(0, 1))}_{=0}] \\
 &= E[Y^2] \\
 &= \text{Var}(Y) \quad (\because E[Y] = 0)
 \end{aligned}$$

$$\text{Var}(Y) = \text{Var}(X_1) + \text{Var}(N(0, \sigma^2)) = 2 \cdot \sigma^2$$

$$\begin{aligned}
 \text{Var}(X_2) &= \text{Var}(Y) + \text{Var}(N(0, 1)) \\
 &= 2\sigma^2 + 1
 \end{aligned}$$

$$\therefore \hat{d}_2 = \frac{\sigma^2}{\sigma^2 + 1}$$

Domain-invariant representation learning

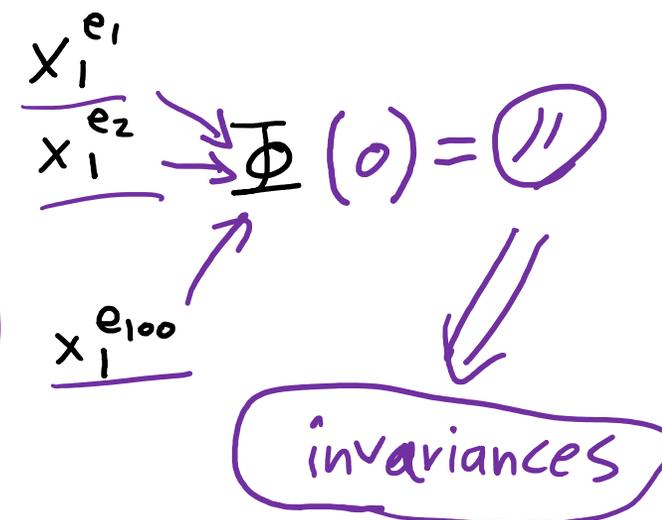
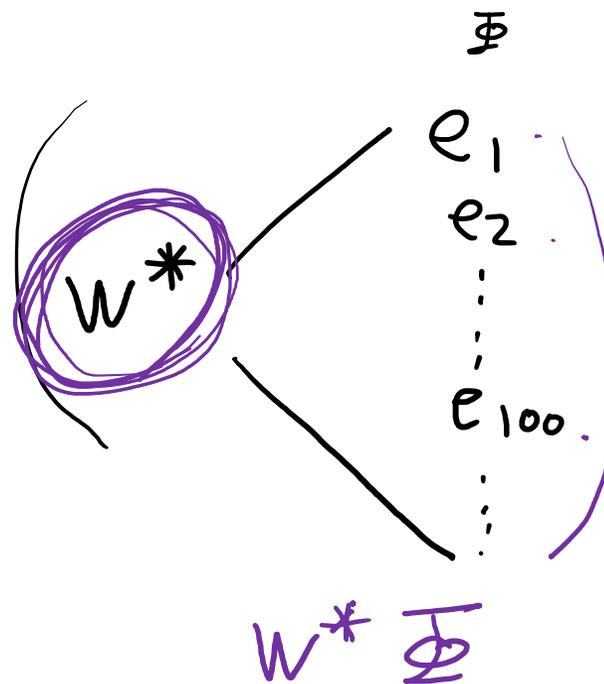
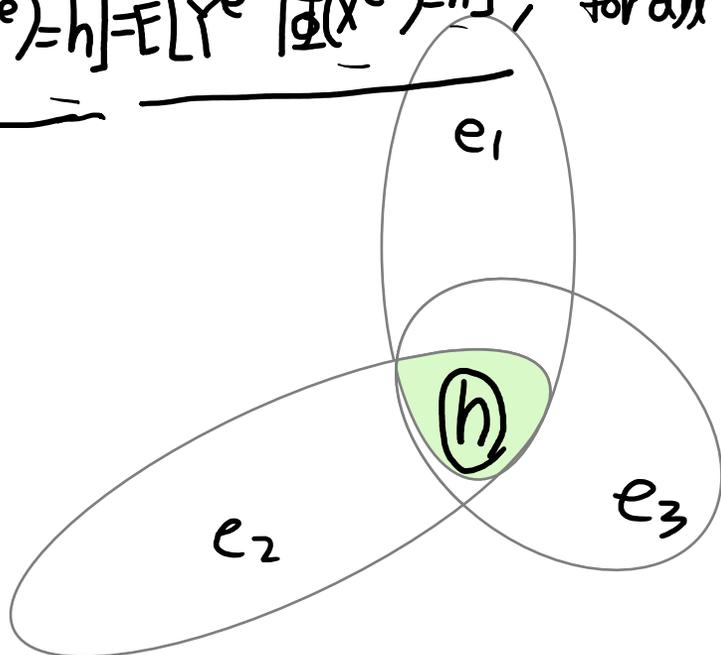
Invariant Risk Minimization

❖ Empirical data로부터 invariances를 얻는 원리

Definition 3. We say that a data representation $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ elicits an invariant predictor $w \circ \Phi$ across environments \mathcal{E} (if there is a classifier $w : \mathcal{H} \rightarrow \mathcal{Y}$ simultaneously optimal for all environments, that is, $w \in \arg \min_{\bar{w} : \mathcal{H} \rightarrow \mathcal{Y}} R^e(\bar{w} \circ \Phi)$ for all $e \in \mathcal{E}$.)

MSE, cross-ent.

$E[Y^e | \Phi(x^e) = h] = E[Y^{e'} | \Phi(x^{e'}) = h]$, for all $e, e' \in \mathcal{E}$



Domain-invariant representation learning

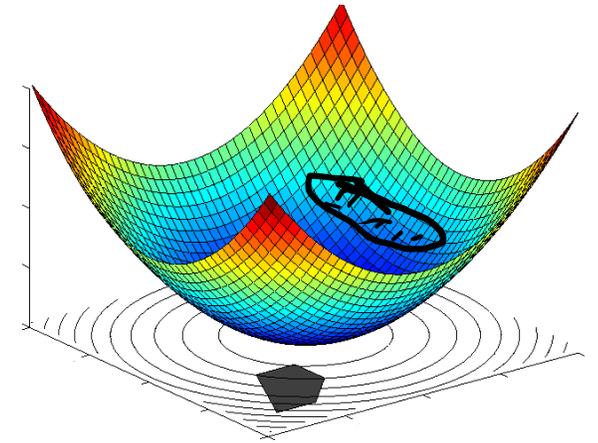
Invariant Risk Minimization

❖ Empirical data로부터 invariances를 얻는 원리

Definition 3. We say that a data representation $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ elicits an invariant predictor $w \circ \Phi$ across environments \mathcal{E} if there is a classifier $w : \mathcal{H} \rightarrow \mathcal{Y}$ simultaneously optimal for all environments, that is, $w \in \arg \min_{\bar{w}: \mathcal{H} \rightarrow \mathcal{Y}} R^e(\bar{w} \circ \Phi)$ for all $e \in \mathcal{E}$.

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

(Inner routine) : 고정된 Φ 로부터 모든 학습 환경에서 loss의 기대값이 가장 작은 w 를 찾고,
 (Outer routine) : 현재 모든 학습 환경에서 최적인 w 를 사용해서 Φ 를 업데이트



Domain-invariant representation learning

Invariant Risk Minimization

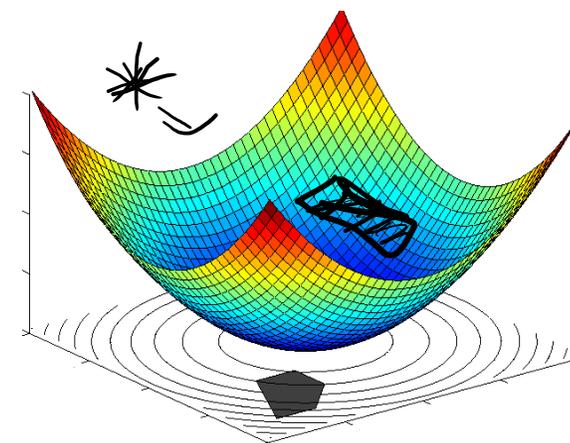
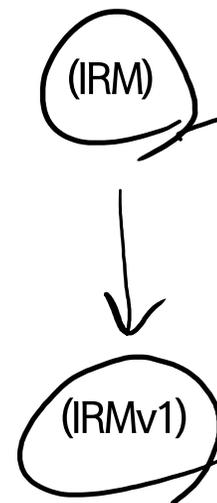
❖ Empirical data로부터 invariances를 얻는 원리

Definition 3. We say that a data representation $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ elicits an invariant predictor $w \circ \Phi$ across environments \mathcal{E} if there is a classifier $w : \mathcal{H} \rightarrow \mathcal{Y}$ simultaneously optimal for all environments, that is, $w \in \arg \min_{\bar{w} : \mathcal{H} \rightarrow \mathcal{Y}} R^e(\bar{w} \circ \Phi)$ for all $e \in \mathcal{E}$.

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

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

ERM과 penalty를 조절하는 hyper-parameter



Domain-invariant representation learning

Invariant Risk Minimization

❖ From (IRM) to (IRMv1)

$$L_{IRM}(\Phi, w) = \sum_{e \in \mathcal{E}_{tr}} (R^e(w \circ \Phi)) + \lambda \cdot \mathbb{D}(w, \Phi, e)$$

Soft regularization term

Classifier w 를 linear로 가정 ✓

Fixed $\Phi, w_{\Phi}^e \in \operatorname{argmin}_{\bar{w}} R^e(\bar{w} \circ \Phi)$: $w_{\Phi}^e = \mathbb{E}_{X^e} [\Phi(X^e)\Phi(X^e)^T]^{-1} \mathbb{E}_{X^e, Y^e} [\Phi(X^e)Y^e]$

$$\mathbb{D}_{dist}(w, \Phi, e) = \|w - w_{\Phi}^e\|^2$$

고정된 Φ 에서의 optimal classifier w_{Φ}^e 와 선택된 w 사이의 거리

$$\min_{\substack{\Phi: X \rightarrow \mathcal{H} \\ w: \mathcal{H} \rightarrow Y}} \sum_{e \in \mathcal{E}_{tr}} R^e(w \circ \Phi) \quad \text{ERM}$$

subject to $w \in \operatorname{argmin}_{\bar{w}: \mathcal{H} \rightarrow Y} R^e(\bar{w} \circ \Phi)$ for all $e \in \mathcal{E}_{tr}$ (IRM)

[Normal equations]

$$\begin{aligned} \frac{\partial}{\partial w} (XW - Y)^T (XW - Y) \\ = 2X^T XW - 2X^T Y = 0. \\ \therefore X^T XW = X^T Y \\ \therefore W = (X^T X)^{-1} \cdot (X^T Y) \end{aligned}$$

Domain-invariant representation learning

Invariant Risk Minimization

$$Y = 1 \cdot X_1 + 0 \cdot X_2$$

Example 1에서 $w = (1, 0)$, data representation $\Phi = \begin{bmatrix} 1 & 0 \\ 0 & c \end{bmatrix}$ 추가

$$Y = \begin{bmatrix} X_1 & X_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & c \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Handwritten notes:

$$\Phi(X^e) = \begin{bmatrix} 1 & 0 \\ 0 & c \end{bmatrix} \begin{bmatrix} X_1^e \\ X_2^e \end{bmatrix} = \begin{bmatrix} X_1^e \\ c \cdot X_2^e \end{bmatrix}$$

$$\Phi(X^e)\Phi(X^e)^T = \begin{bmatrix} X_1^e \\ c \cdot X_2^e \end{bmatrix} \begin{bmatrix} X_1^e & c \cdot X_2^e \end{bmatrix} = \begin{bmatrix} X_1^e X_1^e & c X_1^e X_2^e \\ c X_1^e X_2^e & c^2 X_2^e X_2^e \end{bmatrix}$$

invariance penalty

$$\mathbb{D}_{dist}(w, \Phi, e) = \|w - w_{\Phi}^e\|^2$$

$$w_{\Phi}^e = \mathbb{E}_{X^e}[\Phi(X^e)\Phi(X^e)^T]^{-1} \mathbb{E}_{X^e, Y^e}[\Phi(X^e)Y^e]$$

\mathbb{D}_{dist} 는 $c=0$ 인 지점에서 불연속

$$\mathbb{D}_{lin}(w, \Phi, e) = \|\mathbb{E}_{X^e}[\Phi(X^e)\Phi(X^e)^T]w - \mathbb{E}_{X^e, Y^e}[\Phi(X^e)Y^e]\|^2$$

$$X^T X W = 2 X^T Y$$

Handwritten note: \Downarrow 정규방정식 위반하는 정도

이때, environment e 에 대해서 Least squares predictor $\hat{Y}^e = X_1^e \hat{\alpha}_1 + X_2^e \hat{\alpha}_2$ 를 사용하여 Y 를 예측하면?

- Regress from X_1^e : $\hat{\alpha}_1 = 1, \hat{\alpha}_2 = 0$
- Regress from X_2^e : $\hat{\alpha}_1 = 0, \hat{\alpha}_2 = \sigma(e)^2 / (\sigma(e)^2 + \frac{1}{2})$
- Regress from (X_1^e, X_2^e) : $\hat{\alpha}_1 = 1 / (\sigma(e)^2 + 1), \hat{\alpha}_2 = \sigma(e)^2 / (\sigma(e)^2 + 1)$

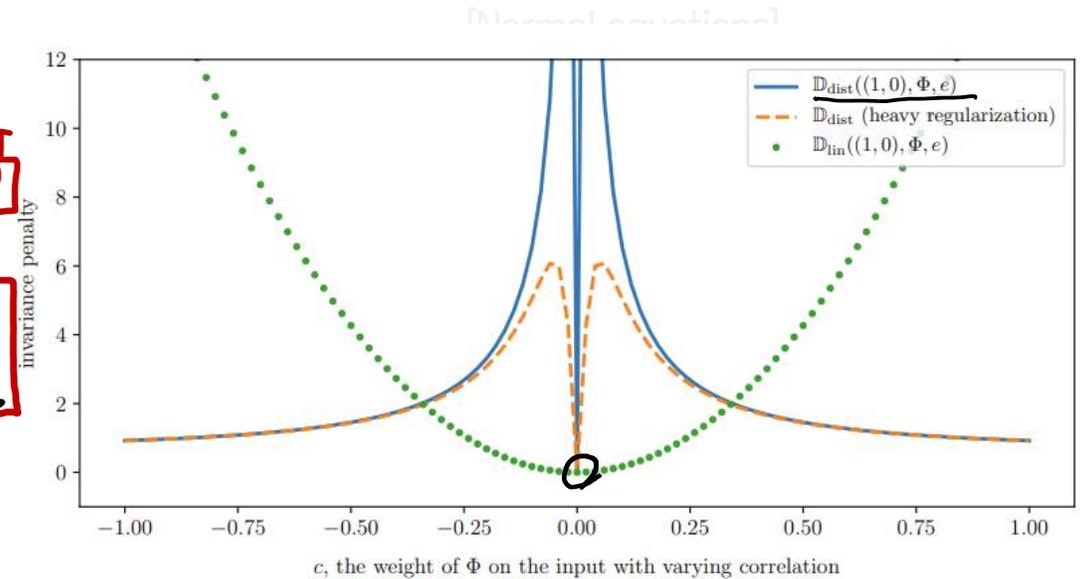


Figure 1: Different measures of invariance lead to different optimization landscapes in our Example 1. The naive approach of measuring the distance between optimal classifiers \mathbb{D}_{dist} leads to a discontinuous penalty (solid blue unregularized, dashed orange regularized). In contrast, the penalty \mathbb{D}_{lin} does not exhibit these problems.



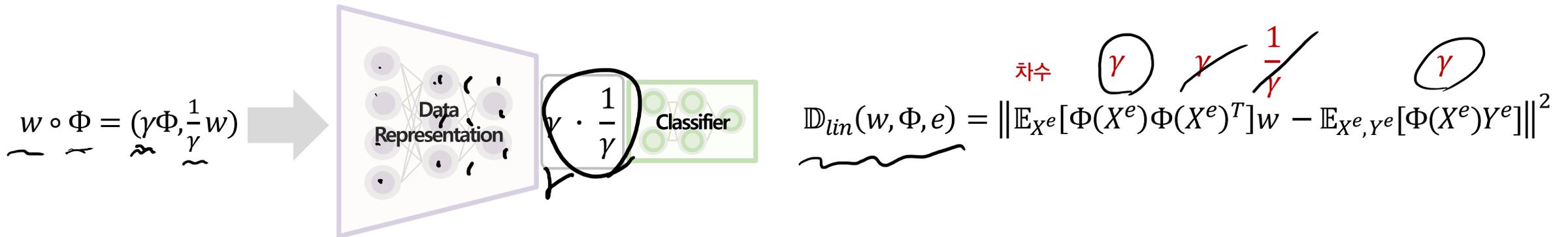
Domain-invariant representation learning

Invariant Risk Minimization

$$L_{IRM}(\Phi, w) = \sum_{e \in \mathcal{E}_{tr}} (R^e(w \circ \Phi) + \lambda \cdot \mathbb{D}_{lin}(w, \Phi, e))$$

❖ Over-parameterized model과 linear classifier w 고정

- Over-parameterized model $\rightarrow \mathbb{D}_{lin} \approx 0$ 가능할 수 있음
- 따라서 classifier w를 고정한 뒤, data representation Φ 만 활용하여 invariant predictor 구축



Over-parameterized model에서는 가중치들이 작은 값을 가져도 표현을 할 수 있음

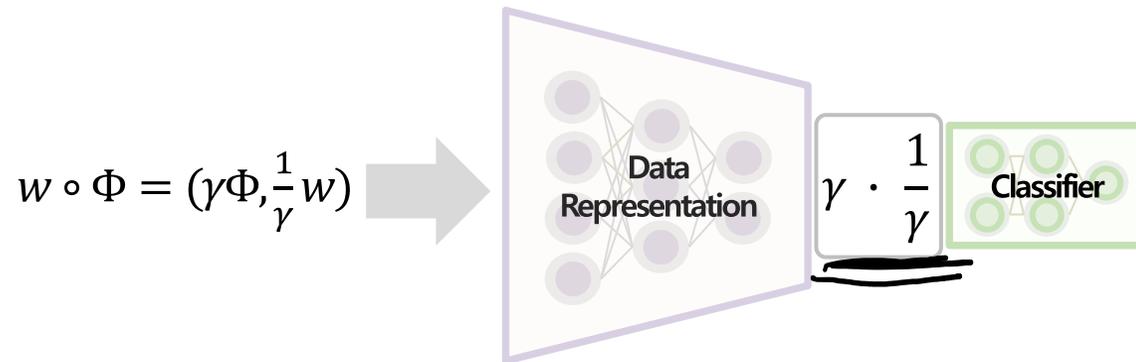
Domain-invariant representation learning

Invariant Risk Minimization

$$L_{IRM}(\Phi, w) = \sum_{e \in \mathcal{E}_{tr}} (R^e(w \circ \Phi) + \lambda \cdot \mathbb{D}_{lin}(w, \Phi, e))$$

❖ Over-parameterized model과 linear classifier w 고정

- Over-parameterized model $\rightarrow \mathbb{D}_{lin} \approx 0$ 가능할 수 있음
- 따라서 classifier w 를 고정한 뒤, data representation Φ 만 활용하여 invariant predictor 구축



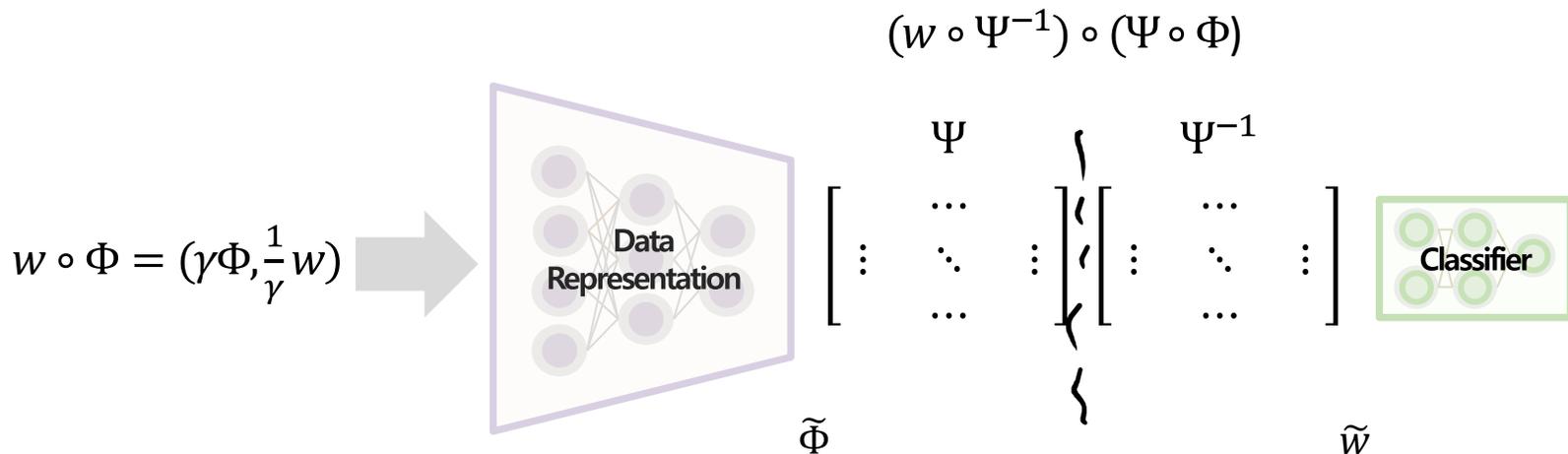
Domain-invariant representation learning

Invariant Risk Minimization

$$L_{IRM}(\Phi, w) = \sum_{e \in \mathcal{E}_{tr}} (R^e(w \circ \Phi) + \lambda \cdot \mathbb{D}_{lin}(w, \Phi, e))$$

❖ Over-parameterized model과 linear classifier w 고정

- Over-parameterized model $\rightarrow \mathbb{D}_{lin} \approx 0$ 가능할 수 있음
- 따라서 classifier w를 고정한 뒤, data representation Φ 만 활용하여 invariant predictor 구축



1. w 를 고정해서, $\mathbb{D}_{lin}(w, \Phi, e) \approx 0$ 를 조금 개선할 수 있음.

2. 모든 e 에 항상 '고정' w.

고정된 \tilde{w} 가 모든 environments에 대해서 optimal classifier가 될 수 있는 $\tilde{\Phi}$ 를 찾음으로써 invariant predictor 구축

$$L_{IRM, w=\tilde{w}}(\Phi) = \sum_{e \in \mathcal{E}_{tr}} (R^e(\tilde{w} \circ \Phi) + \lambda \cdot \mathbb{D}_{lin}(\tilde{w}, \Phi, e))$$



Domain-invariant representation learning

Invariant Risk Minimization

❖ \tilde{w} 를 scalar값으로 고정시켜도 invariance를 학습할 수 있음

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 $v = \Phi^T w$, 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 $v^T \nabla R^e(v) = 0$ 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 $\nabla R^e(v)$.

$$w^* \in \underset{\bar{w} \in \mathbb{R}}{\text{argmin}} R^e(\bar{w} \circ \Phi), \forall e$$

$$\frac{\partial}{\partial w} R^e(\Phi^T w) = \Phi \nabla R^e(\Phi^T w) = 0$$

$$w^T \Phi \nabla R^e(v) = w^T \cdot 0 = \vec{0}$$

$$\therefore \underline{v^T \nabla R^e(v)} = 0$$

$w \rightarrow 1.0$
 $\Phi \rightarrow \text{non-linear}$

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

Domain-invariant representation learning

Invariant Risk Minimization

❖ 훈련 도메인에서 학습한 invariance를 모든 도메인으로 확장

- 필요 조건 1: 훈련 도메인이 충분한 다양성을 가지고 있어야 함
- 필요 조건 2: 다양한 도메인이 invariance를 가지고 있어야 함

Assumption 8. A set of training environments \mathcal{E}_{tr} lie in linear general position of degree r if $|\mathcal{E}_{tr}| > d - r + \frac{d}{r}$ for some $r \in \mathbb{N}$, and for all non-zero $x \in \mathbb{R}^d$.

훈련 도메인 수

$$\dim \left(\text{span} \left(\left\{ \mathbb{E}_{X^e} [X^e X^{e\top}] x - \mathbb{E}_{X^e, \epsilon^e} [X^e \epsilon^e] \right\}_{e \in \mathcal{E}_{tr}} \right) \right) > d - r.$$

Span

훈련 도메인으로부터 생성된 벡터 공간의 차원

Conclusion

❖ Summary

- Empirical risk minimization은 학습과 테스트 데이터가 같은 분포에서 i.i.d. 샘플링 되었다는 가정이 필요함
- 이러한 가정이 깨져서 발생하는 두 데이터셋 사이의 distribution shift 상황에서 ERM은 좋지 못한 성능을 보임
- Distribution shift 상황에서도 테스트 데이터에 대한 성능을 개선할 수 있는 domain generalization 소개
 - 서로 분포가 다른 다양한 도메인에서 invariance correlation을 학습하는 방법론
 - Invariance를 복잡한 최적화 상황에서 그래디언트 기반 relaxed version 제안

고맙습니다

