

# Controllable Invariance through Adversarial Feature Learning

Qizhe Xie, Zihang Dai, Yulun Du, Eduard Hovy, Graham  
Neubig

Carnegie Mellon University  
Language Technologies Institute

NIPS 2017

## Introduction

Introduction

## Adversarial Invariant Feature Learning

Framework

Theoretical analysis

## Experiments

Experiments: Fairness Classifications

Experiments: Multi-lingual Machine Translation

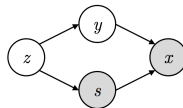
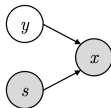
Experiments: Image Classification

- ▶ Representations with invariance properties are often desired
  - ▶ Spatial invariance: CNN
  - ▶ Temporal invariance: RNN
- ▶ This work: a generic framework to induce invariance to a specific factor/attribute of data
  - ▶ Image classifications: classifying people's identities invariant to lighting conditions
  - ▶ Multi-lingual machine translation (fr-en, de-en): translation invariant to source language for sentences with the same meaning
  - ▶ Fairness classifications: predicting credit and saving conditions invariant to the age, gender and race of a person

## Task:

- ▶ Given input  $x$  (images, sentences or features), attribute  $s$  (can be discrete, continuous or structured) of  $x$
- ▶ Predict target  $y$
- ▶ Prior belief: Prediction **should** be invariant to  $s$
- ▶ e.g., predicting identities of a person in a image.  $s$  is the lighting condition
- ▶ Two possible data generation processes:

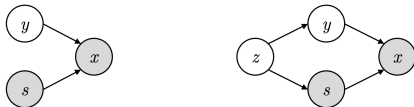
Possible generation process



# Discriminative model

- ▶  $y$  and  $s$  are not independent given  $x$  although they can be marginally independent (Explaining-away)
- ▶  $p(y | x, s)$  is more accurate than  $p(y | x)$ , i.e., knowing  $s$  helps in inferring  $y$ .
  - ▶ “brighten” the representation if it knows the original picture is dark

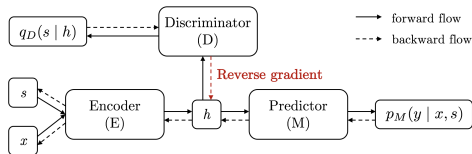
Possible generation process



- ▶ Encoder  $E$ : obtain the invariant representation  $h = E(x, s)$ . ( $s$  is used as the input of the encoder)
- ▶ Predictor  $M$ : Outputs  $q_M(y | h)$  (predict  $y$  based on  $h$ )

## Enforcing Invariance

- ▶  $h$  is invariant to  $s$  means that  $\exists f : f(h) = s$
- ▶ Employ a Discriminator  $D$  to model  $f$ : Outputs  $q_D(s | h)$  (predict  $s$  based on  $h$ )
- ▶ An adversarial game to enforce invariance:
  - ▶ Discriminator tries to detect  $s$  from the representation
  - ▶ Encoder learns to conceal it



## Two objective

- ▶ Standard MLE loss:  $\min_{E, M} -\log q_M(y | h = E(x, s))$
- ▶ Adversarial loss to ensure invariance:  $\min_E \max_D \gamma \log q_D(s | h = E(x, s))$

- ▶ Overall objective:

$$\min_{E, M} \max_D J(E, M, D)$$

where  $J(E, M, D)$  is

$$\mathbb{E}_{x, s, y \sim p(x, s, y)} [\gamma \log q_D(s | h = E(x, s)) - \log q_M(y | h = E(x, s))]$$

- ▶ Definition:  $\tilde{p}(h, s, y) = \int_x p(x, s, y) p_E(h | x, s) dx$
- ▶ Claim 1: Given an encoder, the optimal discriminator and optimal predictor:
  - ▶  $q_D^*(s | h) = \tilde{p}(s | h)$  and  $q_M^*(y | h) = \tilde{p}(y | h)$
  - ▶ Note that  $q_D$  and  $q_M$  are functions of  $E$
- ▶ Claim 2: The optimal encoder is defined by:

$$E^* = \arg \min_E J(E) = \underbrace{-\gamma H(\tilde{q}_E(s | h))}_{\text{Red}} + \underbrace{H(\tilde{q}_E(y | h))}_{\text{Green}}$$

- **[Red]** maximizing the uncertainty of inferring  $s$  based on  $h$
- **[Green]** increasing the certainty of predicting  $y$  based on  $h$

- ▶ The equilibrium of the minimax game is defined by  $\min_E -\gamma H(\tilde{q}(s | h)) + H(\tilde{q}(y | h))$
- ▶ Win-win equilibrium:
  - ▶  $s$  and  $y$  are marginally independent
  - ▶ Two entropy terms reach the optimum at the same time
  - ▶ e.g., removing the lighting conditions in image classifications results in better generalization
- ▶ Competing equilibrium:
  - ▶  $s$  and  $y$  are NOT marginally independent
  - ▶ The optimal of the two entropies cannot be achieved simultaneously
  - ▶ Filtering out  $s$  from  $h$  does harm the prediction of  $y$
  - ▶ e.g., removing bias in fairness classifications hurts the overall performance



# Experiments: Fairness Classifications

- ▶ Task: Predict savings, credit and health condition based on features of a person.  $s$  can be gender or age
- ▶  $E$ ,  $M$ ,  $D$  are all DNN

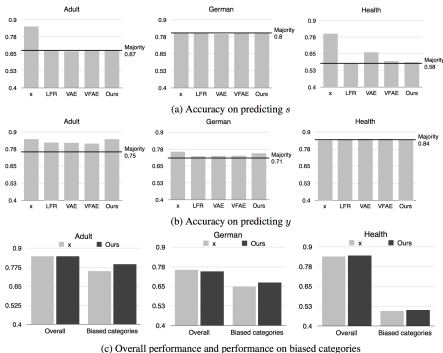


Figure 1: Fair representations should lead to low accuracy on predicting factor  $s$  and high accuracy on predicting  $y$ .

- ▶ Task: Translation from German (de) and French (fr) to English.  $s$  indicates the source language (an one-hot vector)
- ▶  $E$ ,  $M$ ,  $D$  are all LSTM
- ▶ Separate encoders for different languages (Recall that  $h = E(x, s)$ ).
  - ▶ Sharing encoder does not work
  - ▶ DNN based discriminator (even with attention) does not work
  - ▶ Lesson: It is important for  $E$ ,  $M$ ,  $D$  to have enough capacity to achieve the equilibrium

Model	test (fr-en)	test (de-en)
Bilingual Enc-Dec [Bahdanau et al., 2015]	35.2	27.3
Multi-lingual Enc-Dec [Johnson et al., 2016]	35.5	27.7
Our model	<b>36.1</b>	<b>28.1</b>
w.o. discriminator	35.3	27.6
w.o. separate encoders	35.4	27.7

**Table 1:** BLEU score on IWSLT 2015. The ablation study of "w.o. discriminator" shows the improvement is not due to more parameters

# Experiments: Image Classification

- ▶ Task: classifying identities.  $s$  is the lighting condition
- ▶ E, M, D are DNN

Method	Accuracy of classifying factor $s$	Accuracy of classifying target $y$
Logistic regression	0.96	0.78
NN + MMD [Li et al., 2014]	-	0.82
VFAE [Louizos et al., 2016]	<b>0.57</b>	0.85
Ours	<b>0.57</b>	<b>0.89</b>

Table 2: Results on Extended Yale B dataset



Figure 2: t-SNE visualizations of original pictures and learned representations. The original picture is clustered by lighting conditions. The learned representation is clustered by identities