# Mitigating Information Leakage in Image Representations: A Maximum Entropy Approach
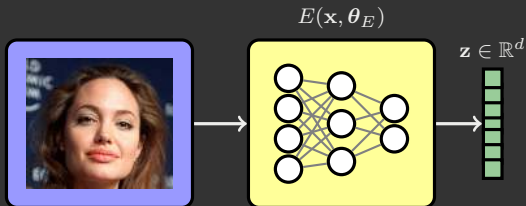
Proteek Roy and Vishnu Boddeti

Michigan State University

CVPR 2019
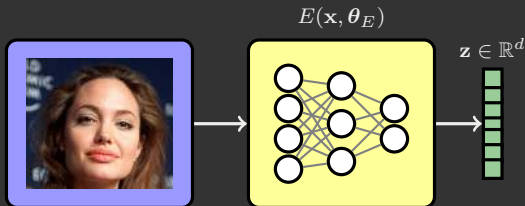
* Deep Embeddings:

* Deep Embeddings:

$$E(\mathbf{x}, \boldsymbol{\theta}_E)$$



$\mathbf{z} \in \mathbb{R}^d$

* Features contain a lot of information
    * basis for generalizing and transferring to other tasks

* Deep Embeddings:



$$E(\mathbf{x}, \boldsymbol{\theta}_E)$$

$$\mathbf{z} \in \mathbb{R}^d$$

* Features contain a lot of information
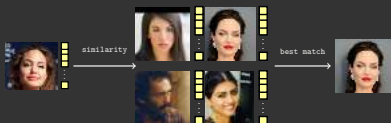    * basis for generalizing and transferring to other tasks

* Applications include:
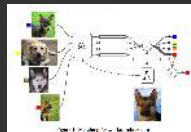


**Figure:** Face Recognition



**Figure:** Image Retrieval

* Features contain a lot of information

```
>>> Representation Learning:  The Dark Side

  * Features contain a lot of information


  * Information may inadvertently be sensitive
```
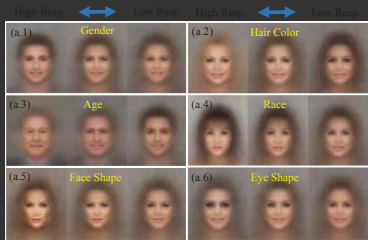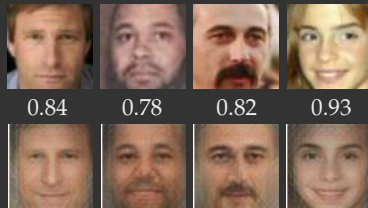
  * Features contain a lot of information

  * Information may inadvertently be sensitive
      * compromise privacy of data owner
      * result in unfair or biased decision systems

* Features contain a lot of information

* Information may inadvertently be sensitive
    * compromise privacy of data owner
    * result in unfair or biased decision systems

* Soft attribute from face features

* Reconstruction from face features



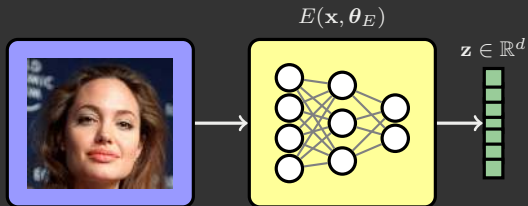Liu et al., ICCV 2015



Mai et al., PAMI 2018

**Mitigating Information Leakage**

Develop representation learning algorithms that can *intentionally* and *permanently* obscure sensitive information while retaining task dependent information.

* Three player zero-sum game between:

$E(\mathbf{x}, \boldsymbol{\theta}_E)$

$\mathbf{z} \in \mathbb{R}^d$

* Three player zero-sum game between:
    * Encoder extracts features $z$

$E(\mathbf{x}, \boldsymbol{\theta}_E)$

$\mathbf{z} \in \mathbb{R}^d$

$T(\boldsymbol{x}, \boldsymbol{\theta}_T)$

$q_T(t|\boldsymbol{z})$

* Three player zero-sum game between:
    * **Encoder** extracts features $z$
    * **Target Predictor** for desired task from features $z$

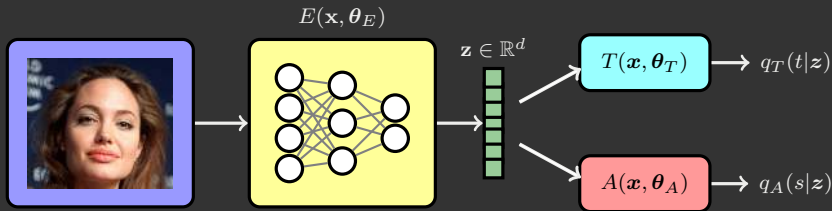$E(\mathbf{x}, \boldsymbol{\theta}_E)$

$\mathbf{z} \in \mathbb{R}^d$

$T(\boldsymbol{x}, \boldsymbol{\theta}_T)$ $\longrightarrow$ $q_T(t|\boldsymbol{z})$

$A(\boldsymbol{x}, \boldsymbol{\theta}_A)$ $\longrightarrow$ $q_A(s|\boldsymbol{z})$

* Three player zero-sum game between:
    * **Encoder** extracts features $\boldsymbol{z}$
    * **Target Predictor** for desired task from features $\boldsymbol{z}$
    * **Adversary** extracts sensitive information from features $\boldsymbol{z}$

$E(\mathbf{x}, \boldsymbol{\theta}_E)$

$\mathbf{z} \in \mathbb{R}^d$

$T(\boldsymbol{x}, \boldsymbol{\theta}_T)$ → $q_T(t|\boldsymbol{z})$

$A(\boldsymbol{x}, \boldsymbol{\theta}_A)$ → $q_A(s|\boldsymbol{z})$
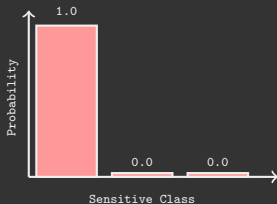
* Three player zero-sum game between:
    * Encoder extracts features $z$
    * Target Predictor for desired task from features $z$
    * Adversary extracts sensitive information from features $z$

* Minimum Likelihood Adversarial Representation Learning:

$$\min_{\boldsymbol{\theta}_E, \boldsymbol{\theta}_T} \max_{\boldsymbol{\theta}_A} \quad \underbrace{J_1(\boldsymbol{\theta}_E, \boldsymbol{\theta}_T)}_{\text{likelihood of predictor}} \quad -\alpha \quad \underbrace{J_2(\boldsymbol{\theta}_E, \boldsymbol{\theta}_A)}_{\text{likelihood of adversary}} \tag{1}$$

* Adversary

* Adversary

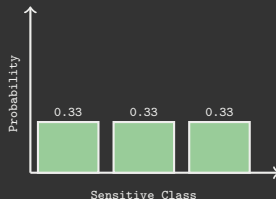* Encoder

* Adversary

* Encoder

* Equillibrium



Limitations:

* Encoder target distribution leaks information !!

* Practice: simultaneous SGD does not reach equilibrium

* Class Imbalance: likelihood biases solution to majority class

Key Idea

Optimize the encoder to maximize entropy of adversary as opposed to minimizing its likelihood.

Key Idea

Optimize the encoder to maximize entropy of adversary as opposed to minimizing its likelihood.

* Adversary

Key Idea

Optimize the encoder to maximize entropy of adversary as
opposed to minimizing its likelihood.

* Adversary

* Encoder

**Key Idea**
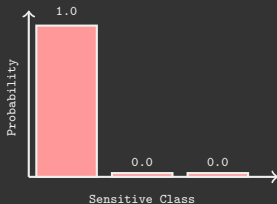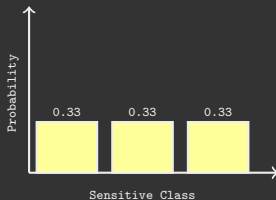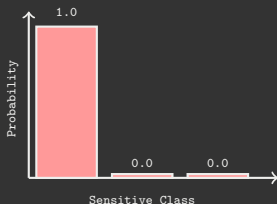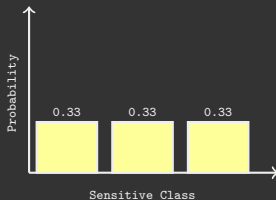
Optimize the encoder to maximize entropy of adversary as opposed to minimizing its likelihood.

* Adversary

* Encoder

* Equilibrium

* Theoretical
    * Three player non-zero sum game
    * At equilibrium, encoder induces uniform distribution in adversary when $s \perp\!\!\!\perp t$
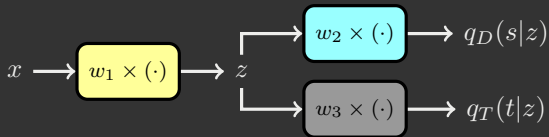    * Obtain conditions for stability of solution around equillibrium through linearization.

* Theoretical
    * Three player non-zero sum game

    * At equilibrium, encoder induces uniform distribution in adversary when $s \perp\!\!\!\perp t$

    * Obtain conditions for stability of solution around equillibrium through linearization.

* Practical
    * Semi-Supervised Mode: encoder does not need sensitive labels

    * Less susceptible to class imbalance than ML-ARL

* Each entity is linear scalar multiplication
* Global solution is $(w_1, w_2, w_3) = (0, 0, 0)$

Minimum Likelihood                    Maximum Entropy

  * UCI Datatset:  Creditworthiness Prediction
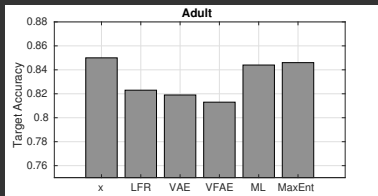
  * UCI Datatset:  Income Prediction

* UCI Datatset:  Creditworthiness Prediction
  Target:  Credit Prediction



* UCI Datatset:  Income Prediction
  Target:  Income Prediction
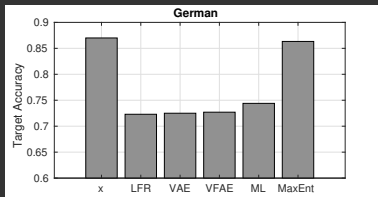
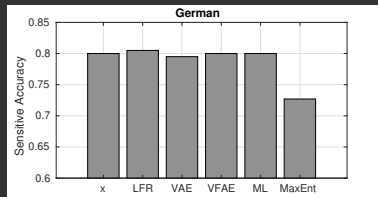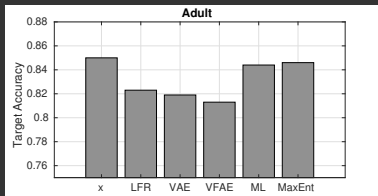* UCI Datatset: Creditworthiness Prediction

    Target: Credit Prediction

    Adversary: Gender Prediction



* UCI Datatset: Income Prediction

    Target: Income Prediction

    Adversary: Gender Prediction

* 38 identities and 5 illumination directions

* Target:   Identity Label

* Sensitive:   Illumination Label

* 38 identities and 5 illumination directions

* Target: Identity Label

* Sensitive: Illumination Label

| Method | $s$ (lighting) | $t$ (identity) |
|---|---|---|
| LR | 96 | 78 |
| NN + MMD (NIPS 2014) | – | 82 |
| VFAE (ICLR 2016) | 57 | 85 |
| ML-ARL (NIPS 2017) | 57 | 89 |
| Maxent-ARL | 40 | 89 |

`>>>` Numerical Experiments: CIFAR-100

- 100 classes categorized into 20 superclasses
- Target: Superclass Label
- Sensitive: Class Label

* 100 classes categorized into 20 superclasses

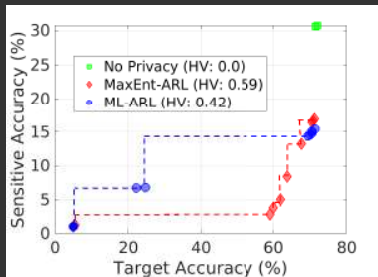* Target:  Superclass Label

* Sensitive:  Class Label



Trade-Off:  Likelihood

* 100 classes categorized into 20 superclasses

* Target:  Superclass Label

* Sensitive:  Class Label



Trade-Off:  Likelihood



Trade-Off:  Entropy

>>> Summary

* A striving step towards explicitly controlling information in learned representations.

* A striving step towards explicitly controlling information in learned representations.

* MaxEnt-ARL: optimize the encoder to maximize entropy of adversary instead of minimizing likelihood.

* A striving step towards explicitly controlling information in learned representations.

* MaxEnt-ARL: optimize the encoder to maximize entropy of adversary instead of minimizing likelihood.

* MaxEnt-ARL enjoys theoretical and practical benefits.

* A striving step towards explicitly controlling information in learned representations.

* MaxEnt-ARL: optimize the encoder to maximize entropy of adversary instead of minimizing likelihood.

* MaxEnt-ARL enjoys theoretical and practical benefits.

Code:

https://github.com/human-analysis/MaxEnt-ARL.git

More Details:   Poster # 175