



UNIVERSITY OF
COPENHAGEN



Building trustworthy ML

The role of label quality and availability

Alexandru Țifrea
Research Scientist,
Google DeepMind

Amartya Sanyal
Tenure Track Assistant Professor,
Department of Computer Science,
University of Copenhagen

Introduction

Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness

Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness



Whether Machine Learning Algorithms have disproportionately worse impact on some groups of people than others

Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness

Privacy



Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness

Privacy



Whether Machine Learning algorithms leak *personal* (training) data

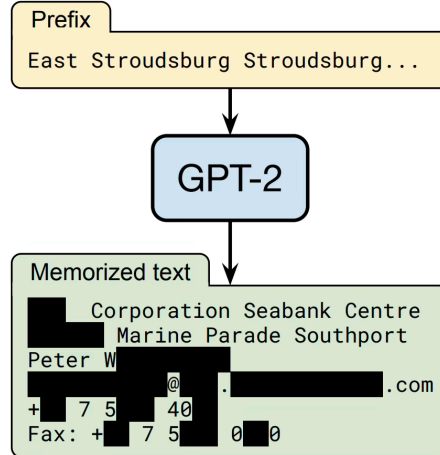
Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness

Privacy

Robustness



Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness

Privacy

Robustness

Whether Machine Learning model can generalise to different data distributions



Aspect 1: Trustworthiness of ML Algorithms

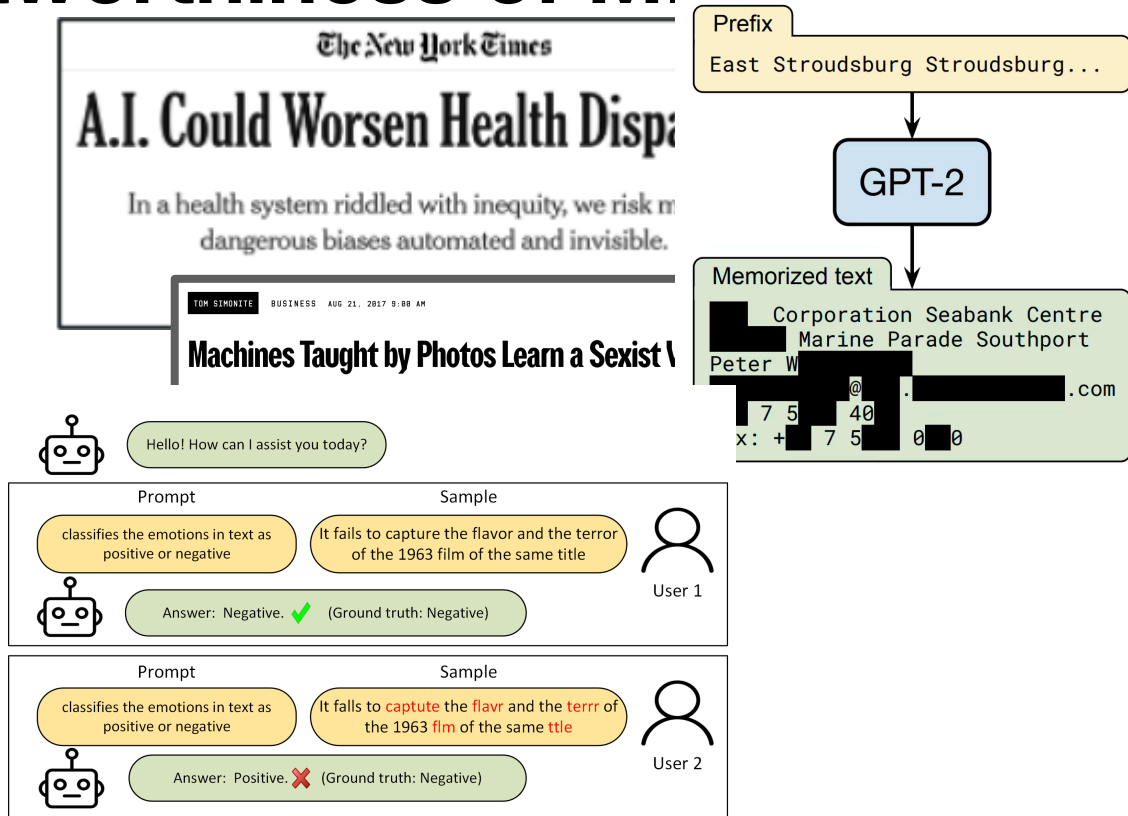
An overloaded term

Fairness

Privacy

Robustness

Whether Machine Learning model can generalise to different data distributions



Aspect 1: Trustworthiness of ML Algorithms

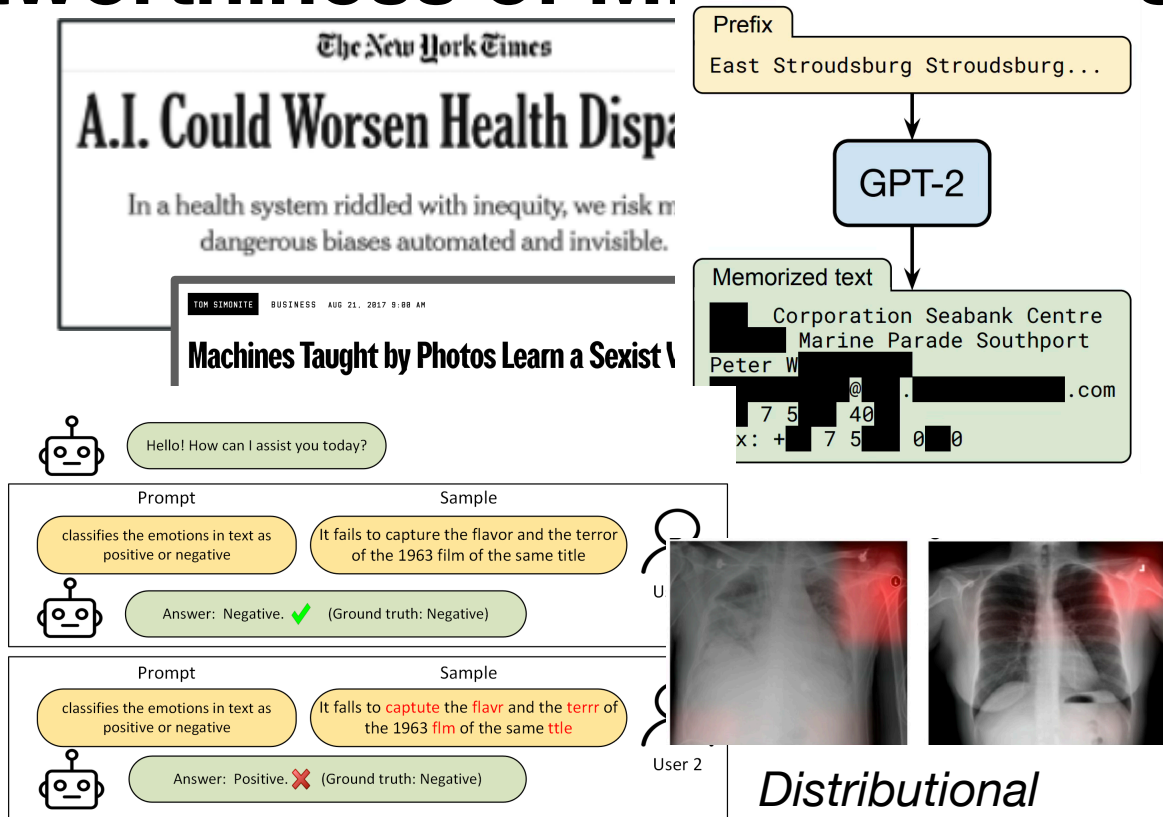
An overloaded term

Fairness

Privacy

Robustness

Whether Machine Learning model can generalise to different data distributions



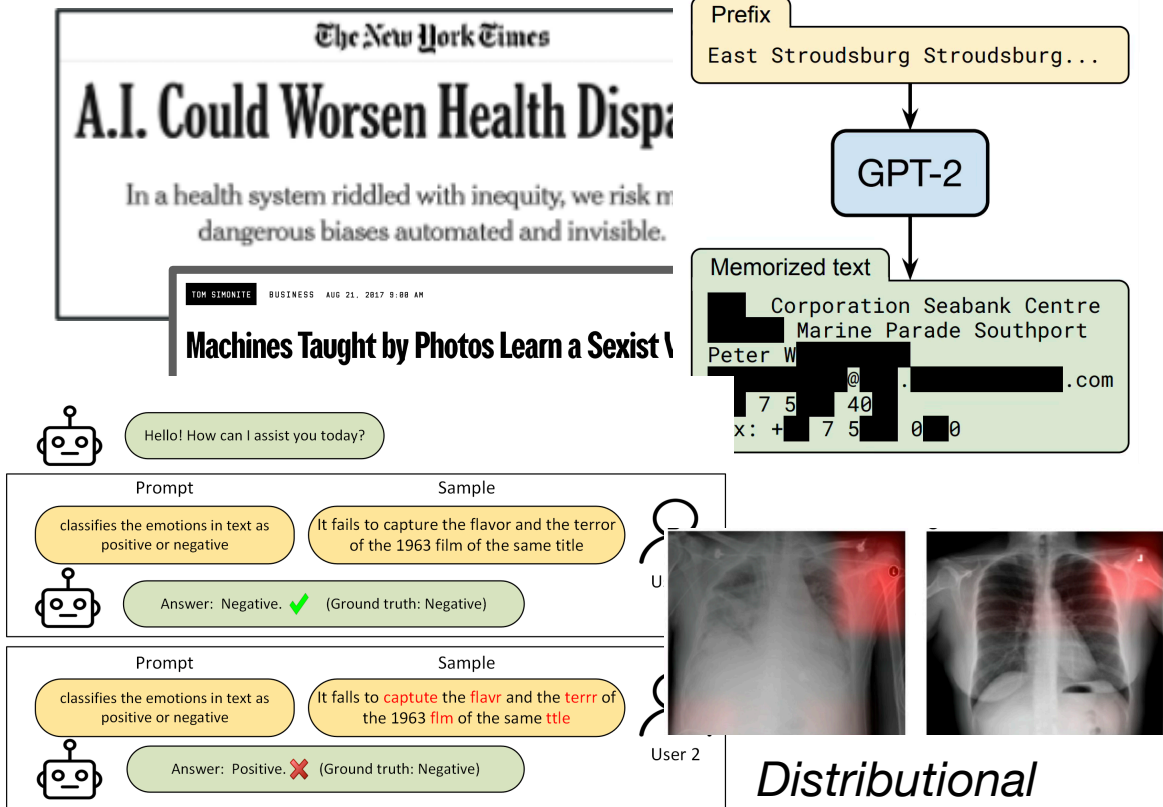
Aspect 1: Trustworthiness of ML Algorithms

An overloaded term

Fairness

Privacy

Robustness



Aspect 2: Label/data quality and availability

Aspect 2: Label/data quality and availability

Two problems with data in ML

Aspect 2: Label/data quality and availability

Two problems with data in ML

- Unlabelled data is significantly more abundant than Labelled data.

Aspect 2: Label/data quality and availability

Two problems with data in ML

- Unlabelled data is significantly more abundant than Labelled data.
- Label noise is ubiquitous in real world data

Aspect 2: Label/data quality and availability

Two problems with data in ML

- Unlabelled data is significantly more abundant than Labelled data.
- Label noise is ubiquitous in real world data

Dataset	Modality	% error
MNIST	image	0.15
CIFAR-10	image	0.54
CIFAR-100	image	5.85
Caltech-256 [†]	image	1.54
ImageNet*	image	5.83
QuickDraw [†]	image	10.12
20news	text	1.09
IMDB	text	2.90
Amazon Reviews [†]	text	3.90
AudioSet	audio	1.35

Aspect 2: Label/data quality and availability

Two problems with data in ML

- Unlabelled data is significantly more abundant than Labelled data.
- Label noise is ubiquitous in real world data

In this tutorial, we will look at

Dataset	Modality	% error
MNIST	image	0.15
CIFAR-10	image	0.54
CIFAR-100	image	5.85
Caltech-256 [†]	image	1.54
ImageNet*	image	5.83
QuickDraw [†]	image	10.12
20news	text	1.09
IMDB	text	2.90
Amazon Reviews [†]	text	3.90
AudioSet	audio	1.35

Aspect 2: Label/data quality and availability

Two problems with data in ML

- Unlabelled data is significantly more abundant than Labelled data.
- Label noise is ubiquitous in real world data

In this tutorial, we will look at

Dataset	Modality	% error
MNIST	image	0.15
CIFAR-10	image	0.54
CIFAR-100	image	5.85
Caltech-256 [†]	image	1.54
ImageNet*	image	5.83
QuickDraw [†]	image	10.12
20news	text	1.09
IMDB	text	2.90
Amazon Reviews [†]	text	3.90
AudioSet	audio	1.35

How **availability and quality of labels** (and data) specifically impact **Fairness, Privacy, and Robustness** of ML Algorithms

Today's Plan

Today's Plan

- Introduction

Today's Plan

- Introduction
- **Fairness in Machine Learning**

Today's Plan

- Introduction
- **Fairness in Machine Learning**
 - Partial group labels
 - No group labels
 - Low-label regime

Today's Plan

- Introduction
 - **Fairness in Machine Learning**
 - **Privacy in Machine Learning**
- Partial group labels
 - No group labels
 - Low-label regime

Today's Plan

- Introduction
- **Fairness in Machine Learning**
- **Privacy in Machine Learning**

- Partial group labels
- No group labels
- Low-label regime

- Privacy and Disparate Impact
- Good data incurs less cost

Today's Plan

- Introduction
- **Fairness in Machine Learning**
- **Privacy in Machine Learning**
- **Robustness in Machine Learning**

- Partial group labels
- No group labels
- Low-label regime

- Privacy and Disparate Impact
- Good data incurs less cost

Today's Plan

- Introduction
- **Fairness in Machine Learning**
- **Privacy in Machine Learning**
- **Robustness in Machine Learning**

- Partial group labels
- No group labels
- Low-label regime

- Privacy and Disparate Impact
- Good data incurs less cost

- Adversarial Robustness
- Distributional Generalisation
- Out-of-distribution detection

Today's Plan

- Introduction
- **Fairness in Machine Learning**
- **Privacy in Machine Learning**
- **Robustness in Machine Learning**
- Outlook and Future Direction

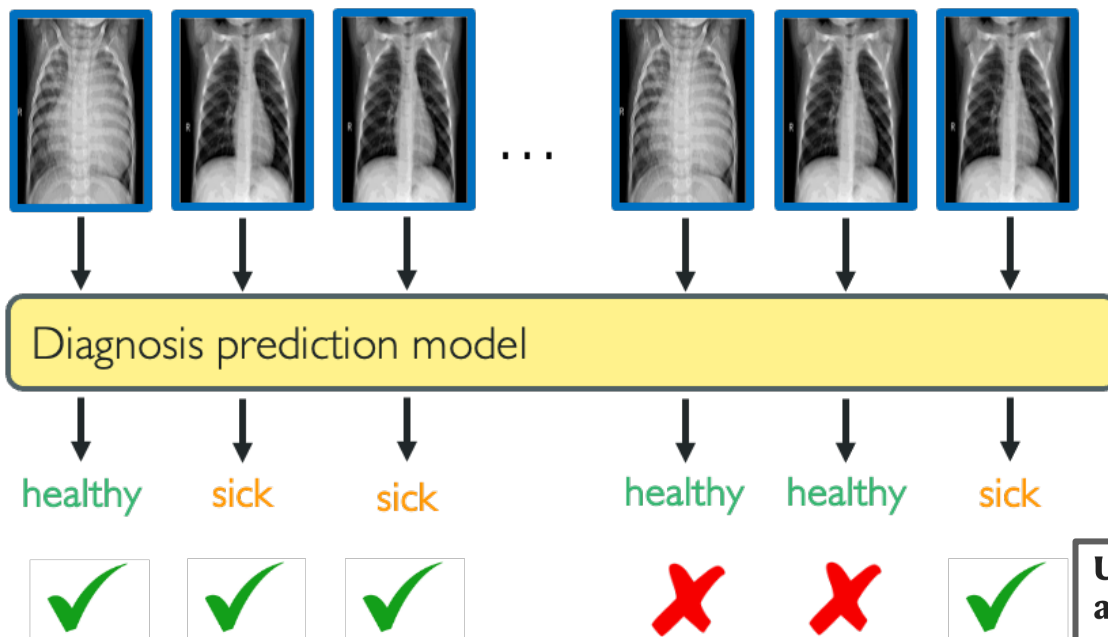
- Partial group labels
- No group labels
- Low-label regime

- Privacy and Disparate Impact
- Good data incurs less cost

- Adversarial Robustness
- Distributional Generalisation
- Out-of-distribution detection

Fairness in Machine Learning

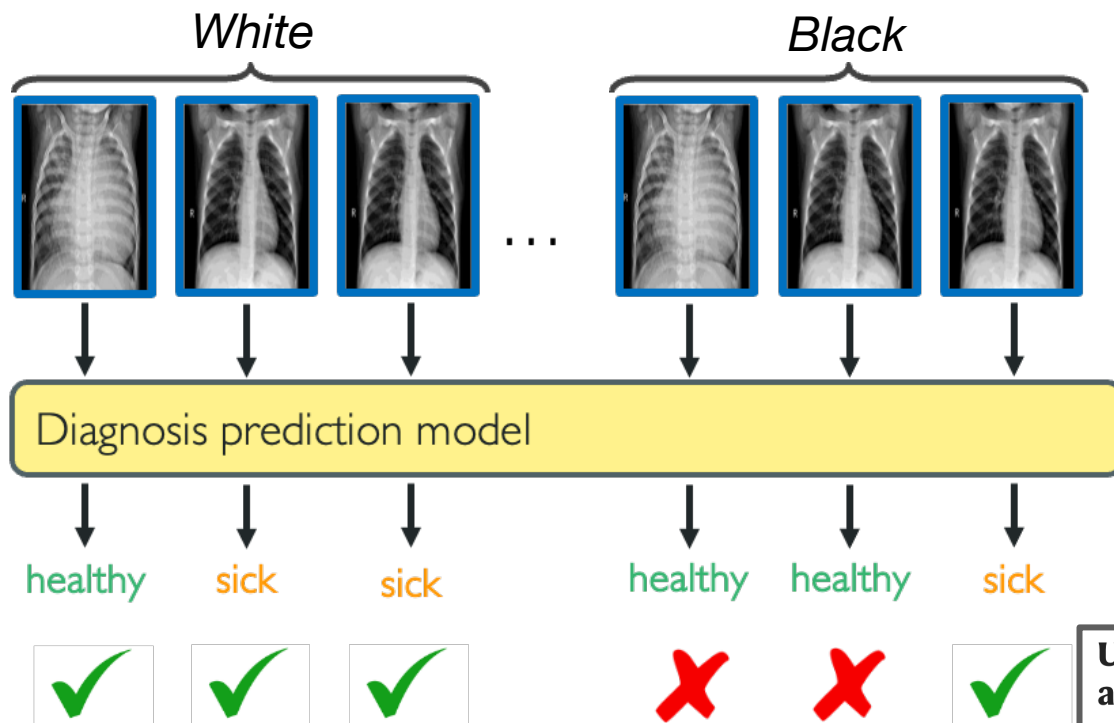
Example of ML model unfairness



Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations

Laleh Seyyed-Kalantari , Haoran Zhang, Matthew B. A. McDermott, Irene Y. Chen & Marzyeh Ghassemi

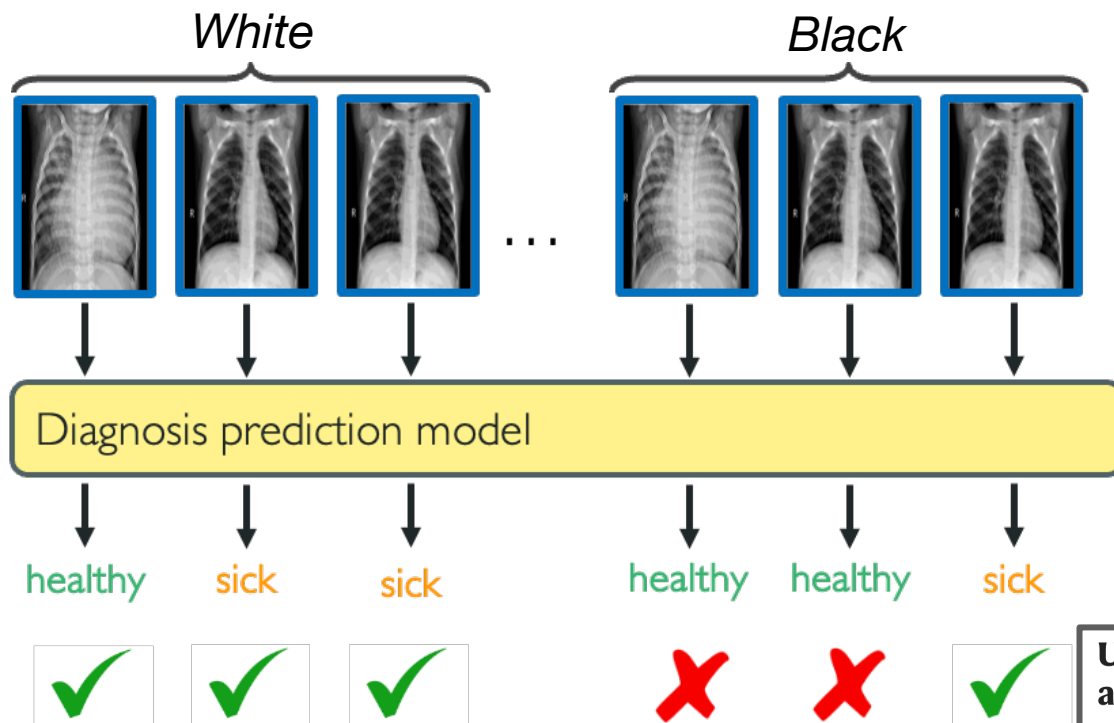
Example of ML model unfairness



Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations

Laleh Seyyed-Kalantari , Haoran Zhang, Matthew B. A. McDermott, Irene Y. Chen & Marzyeh Ghassemi

Example of ML model unfairness



False positive rate:

$FPR = P[\text{predicted healthy} \mid \text{actually sick}]$

$FPR[\text{White}] = 0.16$

$FPR[\text{Black}] = 0.27$

FPR gap = 0.11

**The model is accurate
but not fair!**

Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations

Laleh Seyyed-Kalantari , Haoran Zhang, Matthew B. A. McDermott, Irene Y. Chen & Marzyeh Ghassemi

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Individual fairness: $d_y(\hat{f}(x_1), \hat{f}(x_2)) < C d_x(x_1, x_2)$

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt† Toniann Pitassi‡ Omer Reingold§
Richard Zemel¶

'treating similar individuals similarly'

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Individual fairness: $d_y(\hat{f}(x_1), \hat{f}(x_2)) < C d_x(x_1, x_2)$

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt† Toniann Pitassi‡ Omer Reingold§
Richard Zemel¶

'treating similar individuals similarly'

Group fairness: Three broad categories of fairness notions

Fairness
and
Machine
Learning

o-o

Limitations and Opportunities

Solon Barocas, Moritz Hardt, and Arvind Narayanan

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Individual fairness: $d_y(\hat{f}(x_1), \hat{f}(x_2)) < C d_x(x_1, x_2)$

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt† Toniann Pitassi‡ Omer Reingold§
Richard Zemel¶

'treating similar individuals similarly'

Group fairness: Three broad categories of fairness notions

- Equal acceptance rates
e.g. statistical parity

$$\mathbb{P}(\hat{Y}|A = \text{White}) = \mathbb{P}(\hat{Y}|A = \text{Black})$$

Fairness
and
Machine
Learning

o-o

Limitations and Opportunities

Salon Barocas, Moritz Hardt, and Arvind Narayanan

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Individual fairness: $d_y(\hat{f}(x_1), \hat{f}(x_2)) < C d_x(x_1, x_2)$

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt† Toniann Pitassi‡ Omer Reingold§
Richard Zemel¶

'treating similar individuals similarly'

Group fairness: Three broad categories of fairness notions

- Equal acceptance rates
e.g. statistical parity
- Equal error rates
e.g. Equal Opportunity

$$\mathbb{P}(\hat{Y}|A = \text{White}) = \mathbb{P}(\hat{Y}|A = \text{Black})$$

$$\text{FPR}(A = \text{White}) = \text{FPR}(A = \text{Black})$$

Fairness
and
Machine
Learning

o-o

Limitations and Opportunities

Salon Barocas, Moritz Hardt, and Arvind Narayanan

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Individual fairness: $d_y(\hat{f}(x_1), \hat{f}(x_2)) < C d_x(x_1, x_2)$

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt† Toniann Pitassi‡ Omer Reingold§
Richard Zemel¶

'treating similar individuals similarly'

Group fairness: Three broad categories of fairness notions

- Equal acceptance rates
e.g. statistical parity $\mathbb{P}(\hat{Y}|A = \text{White}) = \mathbb{P}(\hat{Y}|A = \text{Black})$
- Equal error rates
e.g. Equal Opportunity $\text{FPR}(A = \text{White}) = \text{FPR}(A = \text{Black})$
- Equal calibration

Fairness
and
Machine
Learning

o-o

Limitations and Opportunities

Salon Barocas, Moritz Hardt, and Arvind Narayanan

Formal definitions of fairness for prediction

Prediction problem: $\hat{Y} = \hat{f}(X)$ with categorical or continuous labels

Individual fairness: $d_y(\hat{f}(x_1), \hat{f}(x_2)) < C d_x(x_1, x_2)$

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt† Toniann Pitassi‡ Omer Reingold§
Richard Zemel¶

‘treating similar individuals similarly’

Group fairness: Three broad categories of fairness notions

- Equal acceptance rates
e.g. statistical parity $\mathbb{P}(\hat{Y}|A = \text{White}) = \mathbb{P}(\hat{Y}|A = \text{Black})$
- Equal error rates
e.g. Equal Opportunity $\text{FPR}(A = \text{White}) = \text{FPR}(A = \text{Black})$
- Equal calibration

Fairness
and
Machine
Learning

o-o

Limitations and Opportunities

Selon Barocas, Moritz Hardt, and Arvind Narayanan

Remark: Different ML problems (e.g. generative ML) employ similar fairness definitions.

Fairness-error trade-off

State-of-the-art prediction models are often unfair

The New York Times

A.I. Could Worsen Health Disparities

In a health system riddled with inequity, we risk making dangerous biases automated and invisible.

TOM SIMONITE BUSINESS AUG 21, 2017 9:00 AM

Machines Taught by Photos Learn a Sexist View of Women

Algorithms showed a tendency to associate women with shopping and men with shooting.

MIT News
ON CAMPUS AND AROUND THE WORLD

Study reveals why AI models that analyze medical images can be biased

These models, which can predict a patient's race, gender, and age, seem to use those traits as shortcuts when making medical diagnoses.

PRO PUBLICA

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

Fairness-error trade-off

State-of-the-art prediction models are often unfair

The New York Times

A.I. Could Worsen Health Disparities

In a health system riddled with inequity, we risk making dangerous biases automated and invisible.

TOM SIMONITE BUSINESS AUG 21, 2017 9:00 AM

Machines Taught by Photos Learn a Sexist View of Women

Algorithms showed a tendency to associate women with shopping and men with shooting.

MIT News
ON CAMPUS AND AROUND THE WORLD

Study reveals why AI models that analyze medical images can be biased

These models, which can predict a patient's race, gender, and age, seem to use those traits as shortcuts when making medical diagnoses.

PRO PUBLICA

Machine Bias

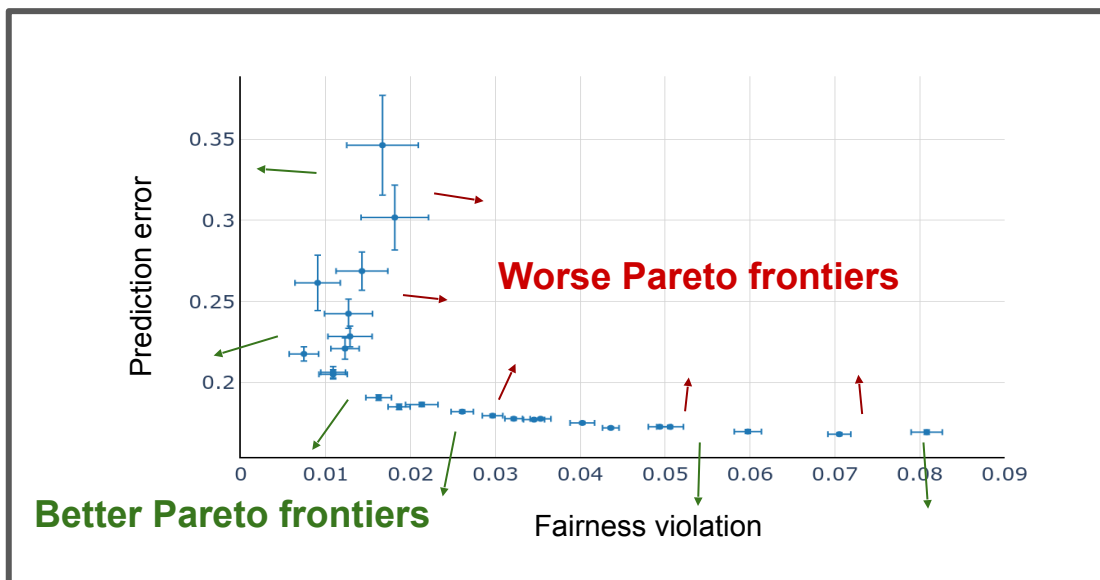
There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

Trivial prediction models (e.g. random guessing) can achieve perfect fairness

e.g. for binary classification and two groups $P(\hat{Y} = 1|A = 0) = P(\hat{Y} = 1|A = 1) = 0.5$

Fairness-error Pareto frontier



need special mitigation algorithms

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \text{ (potentially unfair model)}$$

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

1) **Pre**-processing mitigations

High-level idea: *Change the training data*

Inspired by principle of “Fairness Through Unawareness”

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

1) **Pre-processing** mitigations

High-level idea: *Change the training data*

Inspired by principle of “Fairness Through Unawareness”

Examples:

- feature selection
- fair representation learning
- importance sampling

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \text{ (potentially unfair model)}$$

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

2) **In**-processing mitigations

High-level idea: *Change the training algorithm*

Employ ideas from multi-objective learning

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

2) **In**-processing mitigations

High-level idea: *Change the training algorithm*

Employ ideas from multi-objective learning

e.g. $\arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$ with $\mathcal{D}_{\text{sensitive}} = \{(\mathbf{x}_i, y_i, a_i)\}_{i=1}^m$

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

2) In-processing mitigations

High-level idea: *Change the training algorithm*

Employ ideas from multi-objective learning

e.g. $\arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$ with $\mathcal{D}_{\text{sensitive}} = \{(\mathbf{x}_i, y_i, a_i)\}_{i=1}^m$

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

2) In-processing mitigations

High-level idea: *Change the training algorithm*

Employ ideas from multi-objective learning

e.g. $\arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$ with $\mathcal{D}_{\text{sensitive}} = \{(\mathbf{x}_i, y_i, a_i)\}_{i=1}^m$

unfairness penalty

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

2) In-processing mitigations

High-level idea: *Change the training algorithm*

Employ ideas from multi-objective learning

e.g. $\arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$ with $\mathcal{D}_{\text{sensitive}} = \{(\mathbf{x}_i, y_i, a_i)\}_{i=1}^m$

unfairness penalty

Examples:

- regularized learning
- constrained learning

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \text{ (potentially unfair model)}$$

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

3) **Post**-processing mitigations

High-level idea: *Change the outputs of a pre-trained model*

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

3) **Post**-processing mitigations

High-level idea: *Change the outputs of a pre-trained model*

e.g. group-dependent transformation of outputs:

$$(\hat{Y}, A) \rightarrow T_A(\hat{Y}) \in \{0, 1\}$$

Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \text{ (potentially unfair model)}$$

3) **Post**-processing mitigations

High-level idea: *Change the outputs of a pre-trained model*

e.g. group-dependent transformation of outputs:

$$(\hat{Y}, A) \rightarrow T_A(\hat{Y}) \in \{0, 1\}$$



Fairness mitigation strategies

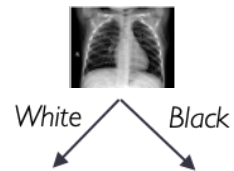
$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

3) **Post**-processing mitigations

High-level idea: *Change the outputs of a pre-trained model*

e.g. group-dependent transformation of outputs:

$$(\hat{Y}, A) \rightarrow T_A(\hat{Y}) \in \{0, 1\}$$



Fairness mitigation strategies

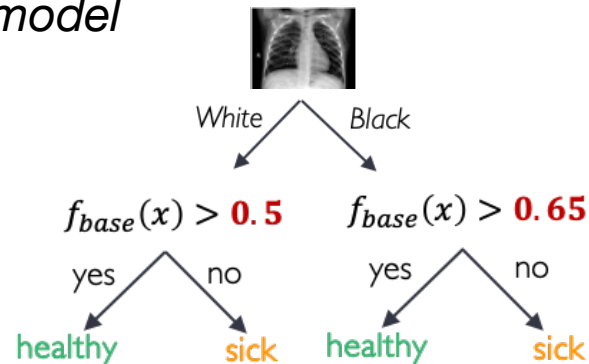
$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

3) **Post**-processing mitigations

High-level idea: *Change the outputs of a pre-trained model*

e.g. group-dependent transformation of outputs:

$$(\hat{Y}, A) \rightarrow T_A(\hat{Y}) \in \{0, 1\}$$



Fairness mitigation strategies

$$\text{OPT}_{\text{base}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}), \quad \mathcal{D}_{\text{pred}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathbb{P}_{XY} \quad (\text{potentially unfair model})$$

3) **Post-processing** mitigations

High-level idea: *Change the outputs of a pre-trained model*

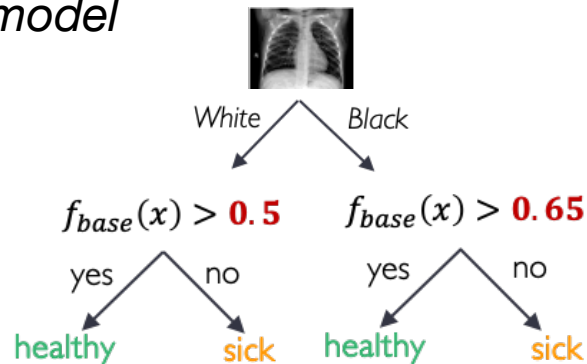
e.g. group-dependent transformation of outputs:

$$(\hat{Y}, A) \rightarrow T_A(\hat{Y}) \in \{0, 1\}$$

Examples:

- group-dependent post-hoc transformations
- group-agnostic transformations

e.g. *fair predictions irrespective of person's willingness to provide sensitive attribute*



Challenges faced by fairness mitigations

Pre-, in-, post-processing mitigations need training data with group labels.

Challenges faced by fairness mitigations

Pre-, in-, post-processing mitigations need training data with group labels.

Issue #1: Group labels are difficult to collect.

e.g. group labels are often sensitive attributes (e.g. gender, ethnicity, age etc)

Challenges faced by fairness mitigations

Pre-, in-, post-processing mitigations need training data with group labels.

Issue #1: Group labels are difficult to collect.

e.g. group labels are often sensitive attributes (e.g. gender, ethnicity, age etc)

What happens when group labels are scarce?

Challenges faced by fairness mitigations

Pre-, in-, post-processing mitigations need training data with group labels.

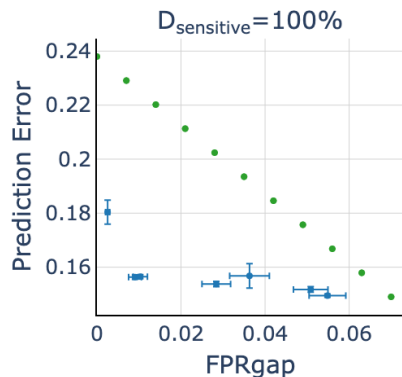
Issue #1: Group labels are difficult to collect.

e.g. group labels are often sensitive attributes (e.g. gender, ethnicity, age etc)

What happens when group labels are scarce?

Naive baseline: “predict according to pre-trained model with probability p , and predict 0 with probability $(1-p)$ ”

In-processing mitigation: state-of-the-art MinDiff method



Dataset: Adult

Y = income; A = gender

Challenges faced by fairness mitigations

Pre-, in-, post-processing mitigations need training data with group labels.

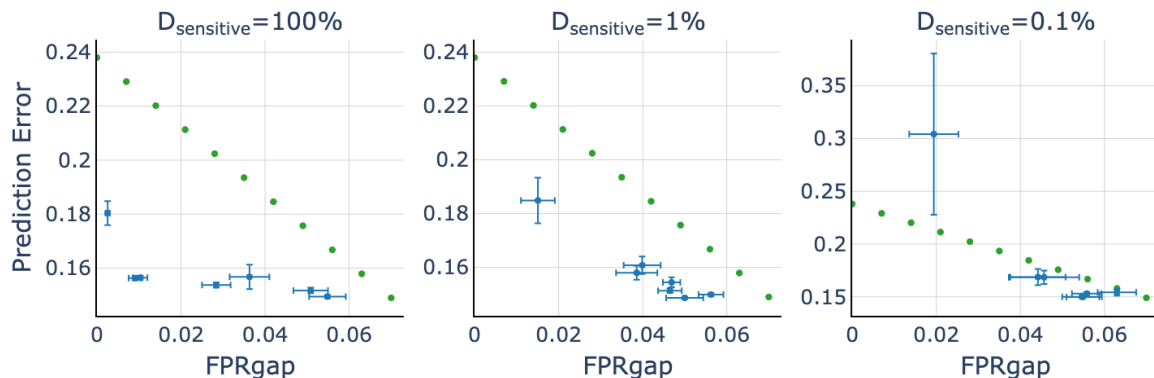
Issue #1: Group labels are difficult to collect.

e.g. group labels are often sensitive attributes (e.g. gender, ethnicity, age etc)

What happens when group labels are scarce?

Naive baseline: “predict according to pre-trained model with probability p , and predict 0 with probability $(1-p)$ ”

In-processing mitigation: state-of-the-art MinDiff method



Dataset: Adult

Y = income; A = gender

Challenges faced by fairness mitigations

Pre-, in-, post-processing mitigations need training data with group labels.

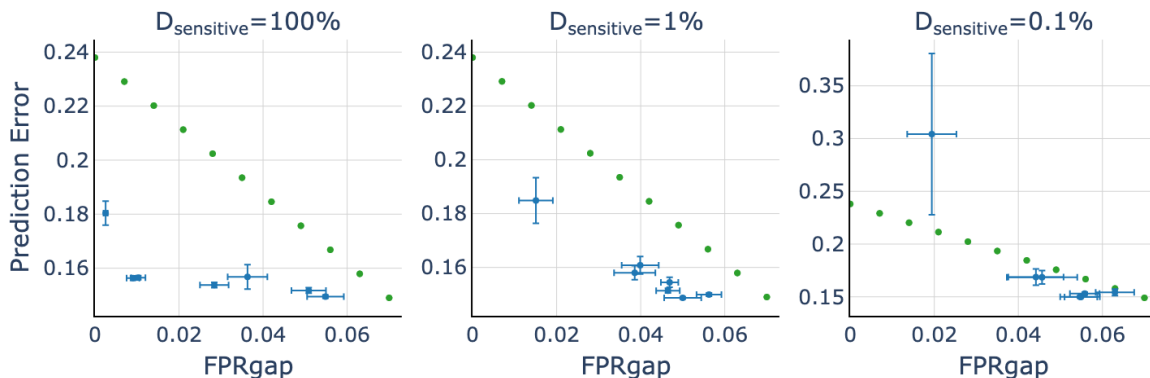
Issue #1: Group labels are difficult to collect.

e.g. group labels are often sensitive attributes (e.g. gender, ethnicity, age etc)

What happens when group labels are scarce?

Naive baseline: “predict according to pre-trained model with probability p , and predict 0 with probability $(1-p)$ ”

In-processing mitigation: state-of-the-art MinDiff method



SOTA method as poor as naive baseline

Dataset: Adult

Y = income; A = gender

Challenges faced by fairness mitigations

Labeled data can be expensive to collect.

Challenges faced by fairness mitigations

Labeled data can be expensive to collect.

Issue #2: Class label scarcity can amplify unfairness.

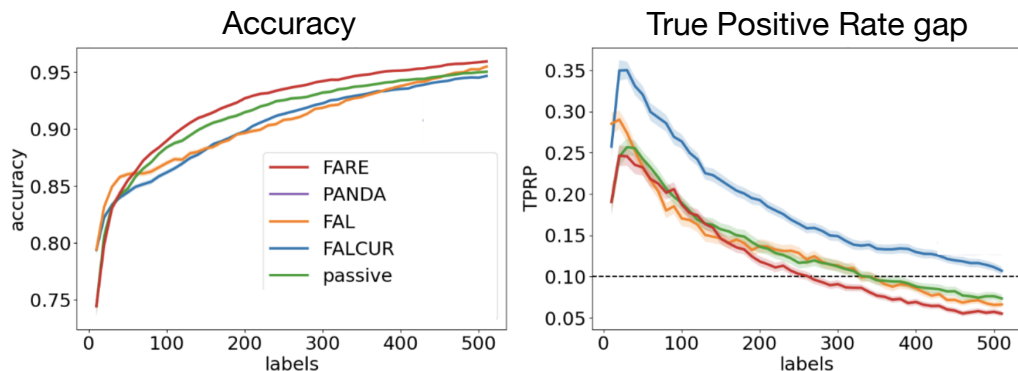
Challenges faced by fairness mitigations

Labeled data can be expensive to collect.

Issue #2: Class label scarcity can amplify unfairness.

What happens in the low-label regime?

e.g. fair active learning strategies



Dataset: Communities & Crime

Y = crime rate; A = ethnicity

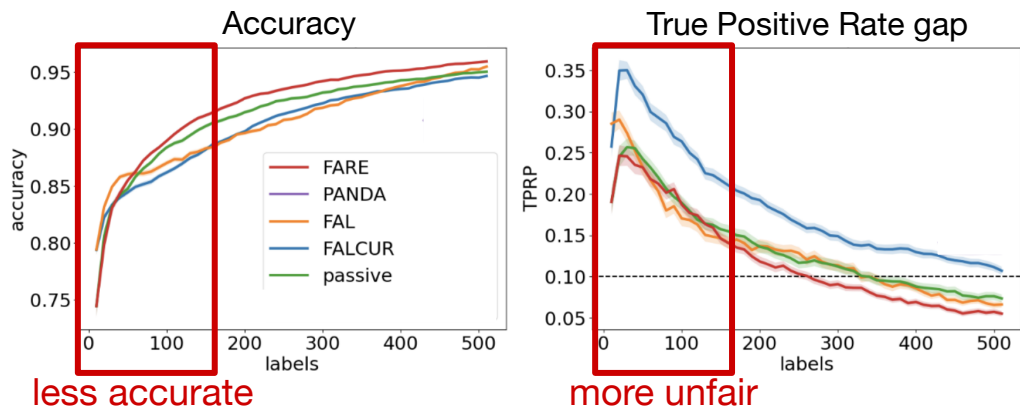
Challenges faced by fairness mitigations

Labeled data can be expensive to collect.

Issue #2: Class label scarcity can amplify unfairness.

What happens in the low-label regime?

e.g. fair active learning strategies



worse accuracy AND fairness in low-label regime

Dataset: Communities & Crime

Y = crime rate; A = ethnicity

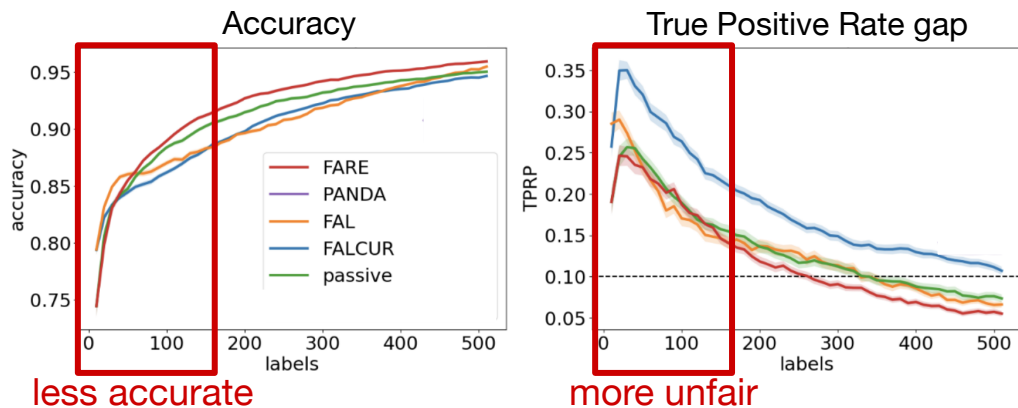
Challenges faced by fairness mitigations

Labeled data can be expensive to collect.

Issue #2: Class label scarcity can amplify unfairness.

What happens in the low-label regime?

e.g. fair active learning strategies



**worse accuracy AND fairness
in low-label regime**

*intersectional fairness amplifies data scarcity
e.g. avoid discriminating against Hispanic females
aged 30-40*

Dataset: Communities & Crime

Y = crime rate; A = ethnicity

Fairness – Outline

Fairness with partial group labels

Fairness with no group labels

Fairness in the low-label regime

Fairness with partial group labels

Problem setting: Fairness with partial group labels

Problem setting: Fairness with partial group labels

$$\mathcal{D}_{\text{pred}} = \{(X_i, Y_i)\}_{i=1}^n \Rightarrow \text{large dataset}$$

covariates X; class labels Y

Problem setting: Fairness with partial group labels

$$\mathcal{D}_{\text{pred}} = \{(X_i, Y_i)\}_{i=1}^n$$

large dataset

covariates X; class labels Y

$$\mathcal{D}_{\text{sensitive}} = \{(X_i, Y_i, A_i)\}_{i=1}^n$$

small dataset

(X, Y) + sensitive attribute A i.e. group label

Problem setting: Fairness with partial group labels

$$\mathcal{D}_{\text{pred}} = \{(X_i, Y_i)\}_{i=1}^n \Rightarrow \text{large dataset}$$

covariates X ; class labels Y

$$\mathcal{D}_{\text{sensitive}} = \{(X_i, Y_i, A_i)\}_{i=1}^n \Rightarrow \text{small dataset}$$

(X, Y) + sensitive attribute A i.e. group label

Case study: In-processing mitigations with partial group labels

Reminder: $\text{OPT}_{\text{IP}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$

Problem setting: Fairness with partial group labels

$$\mathcal{D}_{\text{pred}} = \{(X_i, Y_i)\}_{i=1}^n \Rightarrow \text{large dataset}$$

covariates X ; class labels Y

$$\mathcal{D}_{\text{sensitive}} = \{(X_i, Y_i, A_i)\}_{i=1}^n \Rightarrow \text{small dataset}$$

(X, Y) + sensitive attribute A i.e. group label

Case study: In-processing mitigations with partial group labels

Reminder: $\text{OPT}_{\text{IP}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$

Problem setting: Fairness with partial group labels

$$\mathcal{D}_{\text{pred}} = \{(X_i, Y_i)\}_{i=1}^n \rightarrow \text{large dataset}$$

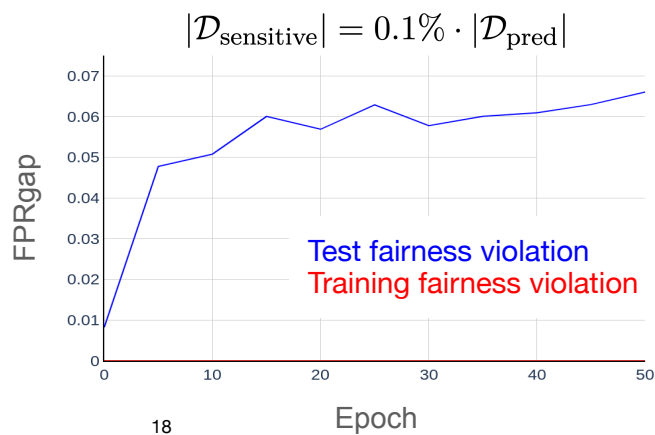
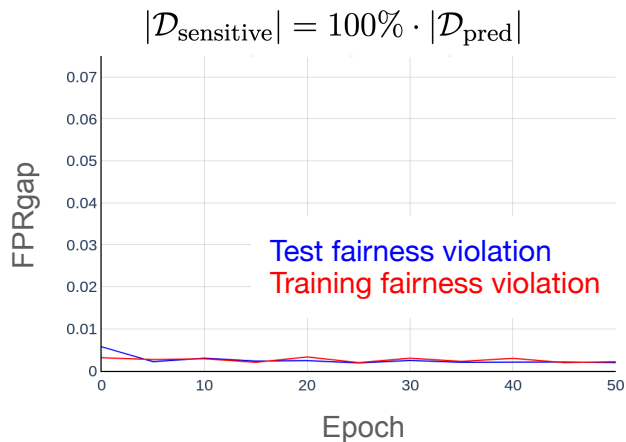
covariates X; class labels Y

$$\mathcal{D}_{\text{sensitive}} = \{(X_i, Y_i, A_i)\}_{i=1}^n \rightarrow \text{small dataset}$$

(X, Y) + sensitive attribute A i.e. group label

Case study: In-processing mitigations with partial group labels

Reminder: $\text{OPT}_{\text{IP}} : \arg \min_f \mathcal{L}_{\text{pred}}(f; \mathcal{D}_{\text{pred}}) + \lambda \mathcal{L}_{\text{fair}}(f; \mathcal{D}_{\text{sensitive}})$



→ overfitting!

How to deal with partial group labels?

High level strategies

1. Use proxy for missing sensitive attributes

1. Make fairness mitigations more sample efficient

How to deal with partial group labels?

High level strategies

1. Use proxy for missing sensitive attributes

1. Make fairness mitigations more sample efficient

Proxy for missing sensitive attributes

Strategies for missing sensitive attributes A

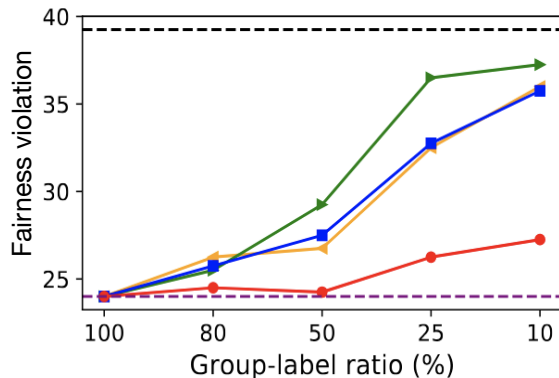
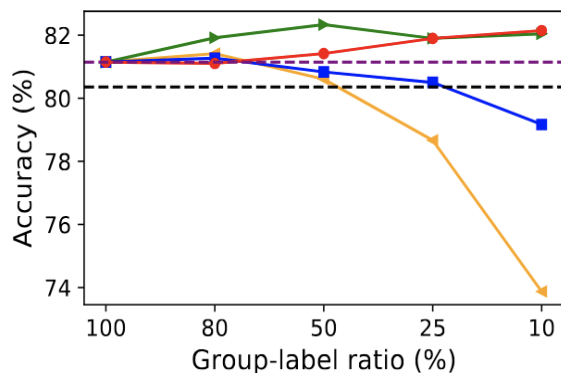
e.g. process data + in-processing fairness mitigation

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University



Dataset: UTKFace

Y = age group; A = ethnicity

Proxy for missing sensitive attributes

Strategies for missing sensitive attributes A

e.g. process data + in-processing fairness mitigation

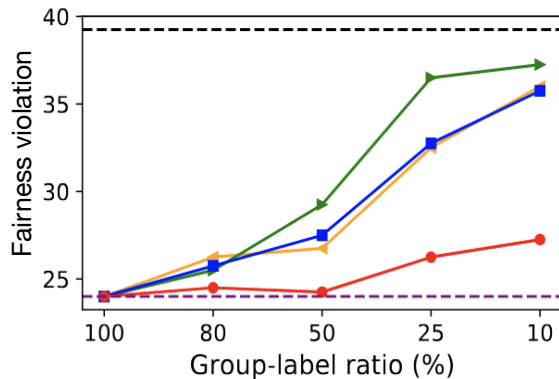
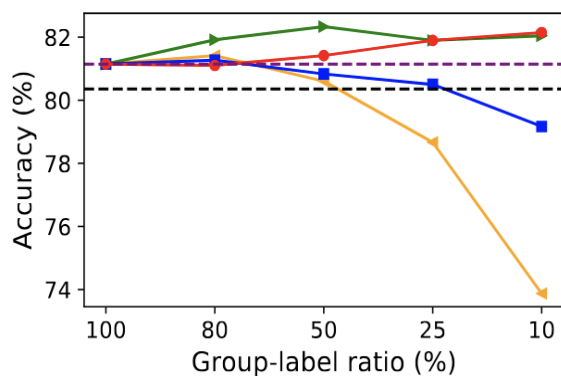
- drop samples with missing A

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University



Dataset: UTKFace

Y = age group; A = ethnicity

Proxy for missing sensitive attributes

Strategies for missing sensitive attributes A

e.g. process data + in-processing fairness mitigation

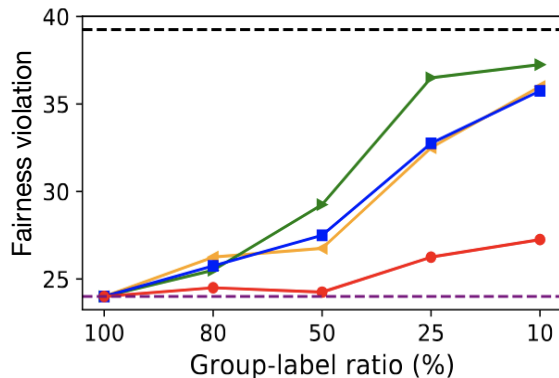
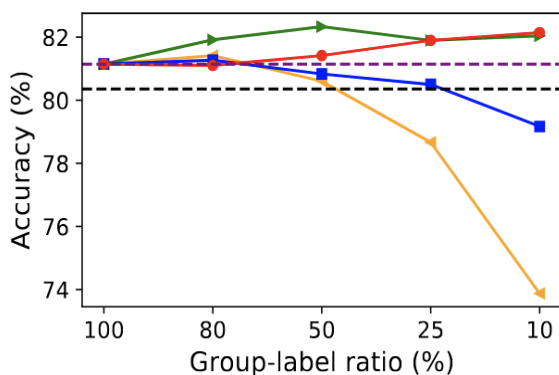
- drop samples with missing A
- impute A uniformly at random

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University



Dataset: UTKFace

Y = age group; A = ethnicity

Proxy for missing sensitive attributes

Strategies for missing sensitive attributes A

e.g. process data + in-processing fairness mitigation

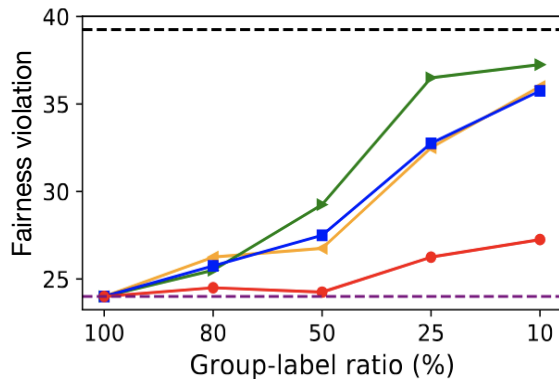
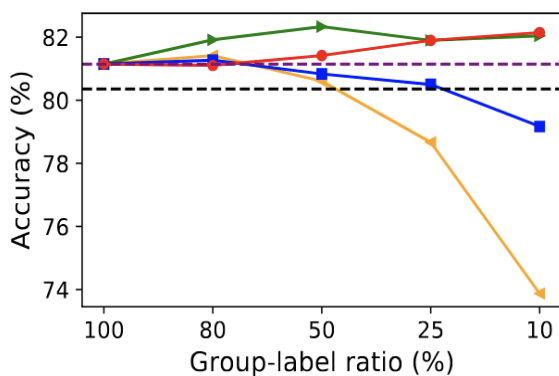
- drop samples with missing A
- impute A uniformly at random
- pseudo-labels from classifier $\hat{f}_A(x)$ trained on $D_{sensitive} = \{(x_i, a_i)\}_{i=1}^n$

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3‡}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University



Dataset: UTKFace

Y = age group; A = ethnicity

Proxy for missing sensitive attributes

Strategies for missing sensitive attributes A

e.g. process data + in-processing fairness mitigation

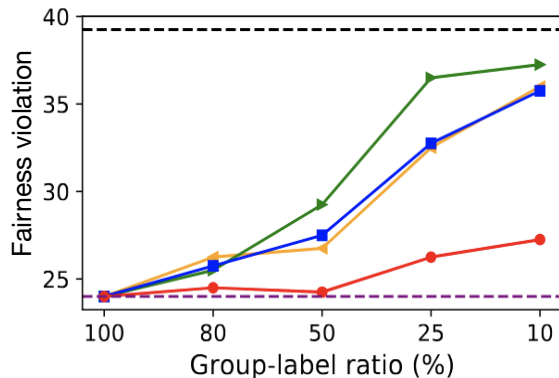
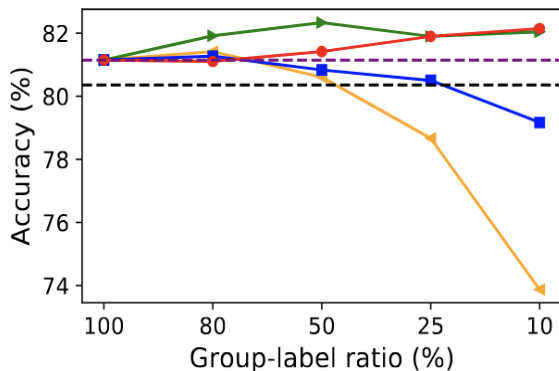
- drop samples with missing A
- impute A uniformly at random
- pseudo-labels from classifier $\hat{f}_A(x)$ trained on $D_{sensitive} = \{(x_i, a_i)\}_{i=1}^n$
- pseudo-labels only on high-confidence samples otherwise random value for A

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University



Dataset: UTKFace

Y = age group; A = ethnicity

Proxy for missing sensitive attributes

Strategies for missing sensitive attributes A

e.g. process data + in-processing fairness mitigation

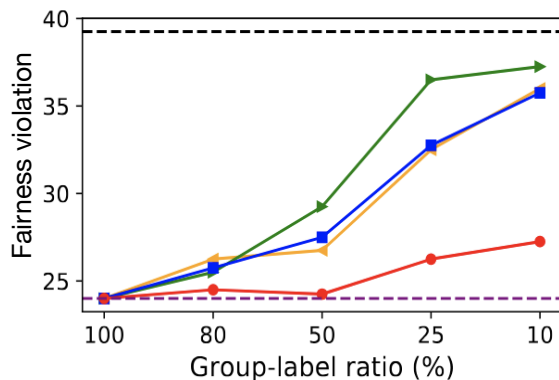
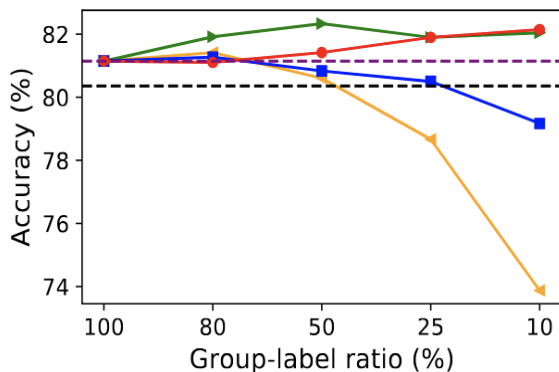
- drop samples with missing A
- impute A uniformly at random
- pseudo-labels from classifier $\hat{f}_A(x)$ trained on $D_{sensitive} = \{(x_i, a_i)\}_{i=1}^n$
- pseudo-labels only on high-confidence samples otherwise random value for A

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University



➔ naive mitigations are suboptimal

Dataset: UTKFace

Y = age group; A = ethnicity

High-confidence group pseudo-labels

Predict missing sensitive attributes A:

$$\hat{a} = \begin{cases} \arg \max \hat{f}_A(x) & \hat{f}_A(x) > \tau \\ \text{draw from } P(A|Y=y) & \text{otherwise} \end{cases}$$

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University

High-confidence group pseudo-labels

Predict missing sensitive attributes A: ?

$$\hat{a} = \begin{cases} \arg \max \hat{f}_A(x) & \hat{f}_A(x) > \tau \\ \text{draw from } P(A|Y=y) & \text{otherwise} \end{cases}$$

Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

³ Interdisciplinary Program in Artificial Intelligence, Seoul National University

High-confidence group pseudo-labels

Predict missing sensitive attributes A: ?

$$\hat{a} = \begin{cases} \arg \max \hat{f}_A(x) & \hat{f}_A(x) > \tau \\ \text{draw from } P(A|Y=y) & \text{otherwise} \end{cases}$$

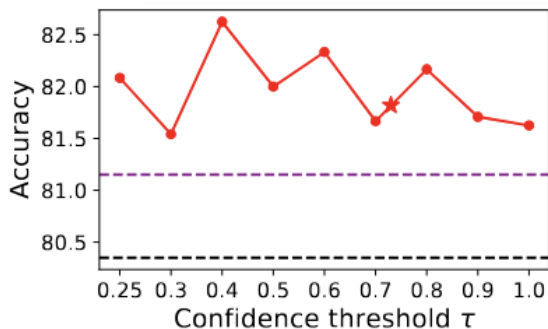
Learning Fair Classifiers with Partially Annotated Group Labels

Sangwon Jung^{1*} Sanghyuk Chun^{2†} Taesup Moon^{1,3†}

¹ Department of ECE/ASRI, Seoul National University ² NAVER AI Lab

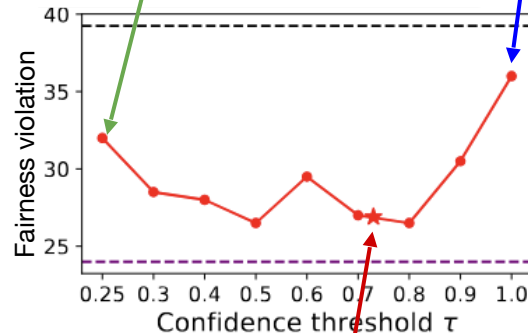
³ Interdisciplinary Program in Artificial Intelligence, Seoul National University

Use validation set
(with group labels)



Suboptimal #1:
always impute
random label

Suboptimal #2:
always impute pseudo-
label from \hat{f}_A



optimal τ^*

Is thresholding confidence an optimal strategy?

Fairness Under Unawareness: Assessing Disparity When Protected Class Is Unobserved

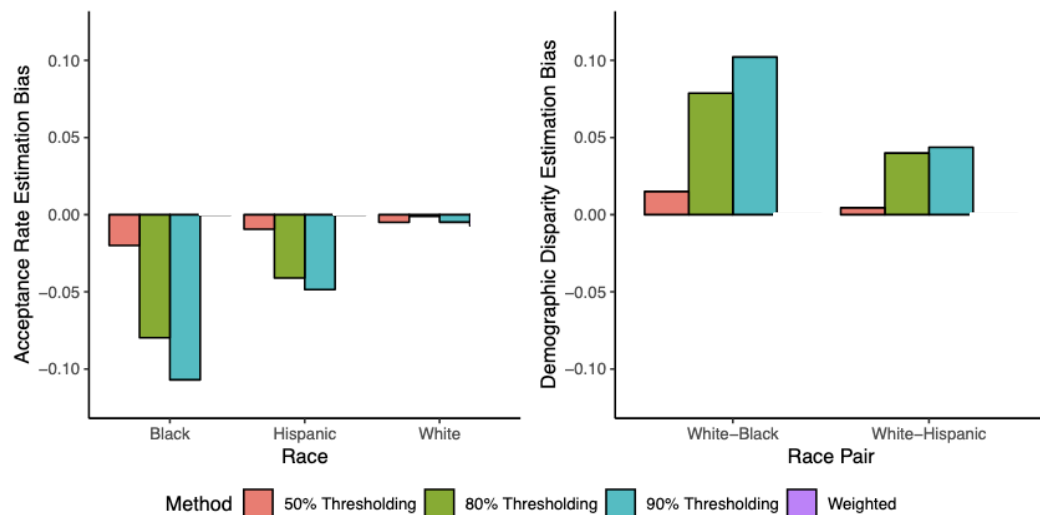
Jiahao Chen
cjiahao@gmail.com

Nathan Kallus
Cornell Tech
New York, New York, USA
kallus@cornell.edu

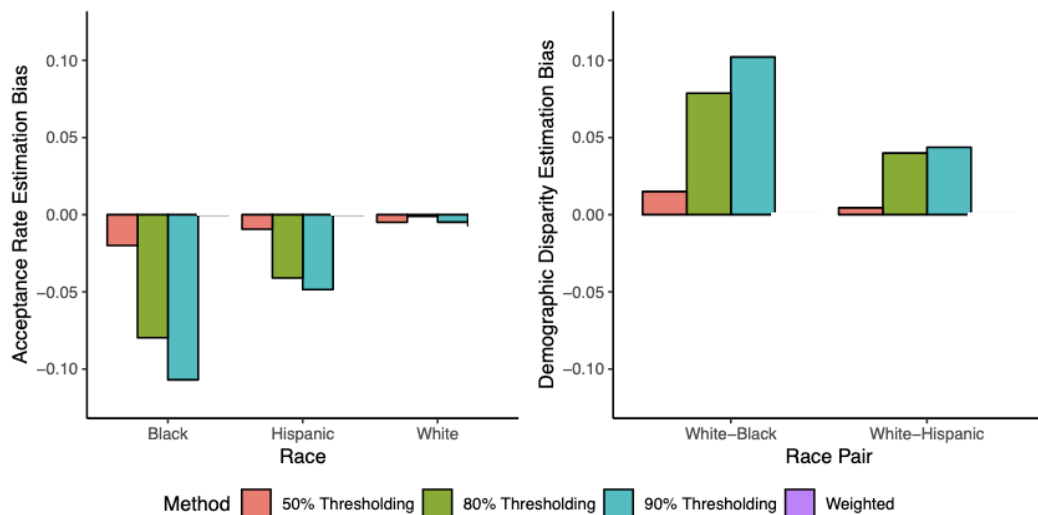
Xiaojie Mao*
Cornell Tech
New York, New York, USA
xm77@cornell.edu

Geoffrey Svacha
svacha@gmail.com

Madeleine Udell
Cornell University
Ithaca, New York, USA
udell@cornell.edu



Is thresholding confidence an optimal strategy?



thresholding confidence can lead to poor fairness estimation

Fairness Under Unawareness: Assessing Disparity When Protected Class Is Unobserved

Jiahao Chen
cjiahao@gmail.com

Nathan Kallus
Cornell Tech
New York, New York, USA
kallus@cornell.edu

Xiaojie Mao*
Cornell Tech
New York, New York, USA
xm77@cornell.edu

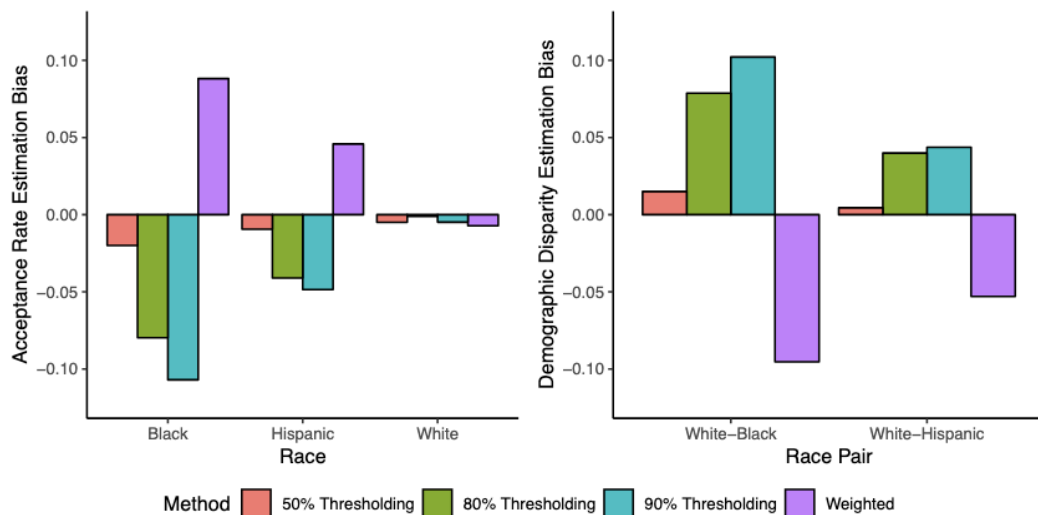
Geoffry Svacha
svacha@gmail.com

Madeleine Udell
Cornell University
Ithaca, New York, USA
udell@cornell.edu

Dataset: HMDA

Y = 'was loan approved?'
A = ethnicity

Is thresholding confidence an optimal strategy?



thresholding confidence can lead to poor fairness estimation

Weighted estimator also fails (but differently)

Fairness Under Unawareness: Assessing Disparity When Protected Class Is Unobserved

Jiahao Chen
cjiahao@gmail.com

Nathan Kallus
Cornell Tech
New York, New York, USA
kallus@cornell.edu

Xiaojie Mao*
Cornell Tech
New York, New York, USA
xm77@cornell.edu

Geoffry Svacha
svacha@gmail.com

Madeleine Udell
Cornell University
Ithaca, New York, USA
udell@cornell.edu

Dataset: HMDA

Y = 'was loan approved?'
A = ethnicity

Summary: Using a proxy group label

Effective at mitigating unfairness

as long as sufficient group-labeled validation data is available

e.g. necessary to select hyperparameters like confidence threshold

Summary: Using a proxy group label

Effective at mitigating unfairness


as long as sufficient group-labeled validation data is available

e.g. necessary to select hyperparameters like confidence threshold

Statistically, often easy to predict the sensitive attribute from little data

but it can have ethical concerns and can amplify/hide biases in the data

Fairer machine learning in the real world: Mitigating discrimination without collecting sensitive data

Michael Veale ¹ and Reuben Binns²

Fairness Under Unawareness:
Assessing Disparity When Protected Class Is Unobserved

Jiahao Chen
cjiahao@gmail.com

Nathan Kallus
Cornell Tech
New York, New York, USA
kallus@cornell.edu

Xiaojie Mao^{*}
Cornell Tech
New York, New York, USA
xm77@cornell.edu

Geoffry Svacha
svacha@gmail.com

Madeleine Udell
Cornell University
Ithaca, New York, USA
udell@cornell.edu

Improving Fairness in Machine Learning Systems:
What Do Industry Practitioners Need?

Kenneth Holstein
Carnegie Mellon University
Pittsburgh, PA
kjholste@cs.cmu.edu

Jennifer Wortman Vaughan
Microsoft Research
New York, NY
jenn@microsoft.com

Hal Daumé III
Microsoft Research &
University of Maryland
New York, NY
me@hal3.name

Miroslav Dudik
Microsoft Research
New York, NY
mdudik@microsoft.com

Hanna Wallach
Microsoft Research
New York, NY
wallach@microsoft.com

How to deal with partial group labels?

High level strategies

1. Use proxy for missing sensitive attributes

- 1. Make fairness mitigations more sample efficient**

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîfrea¹ Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Setup: Equal Opportunity on Adult dataset

Naive baseline: “predict according to f_{base} with probability p , and predict 0 with probability $(1-p)$ ”



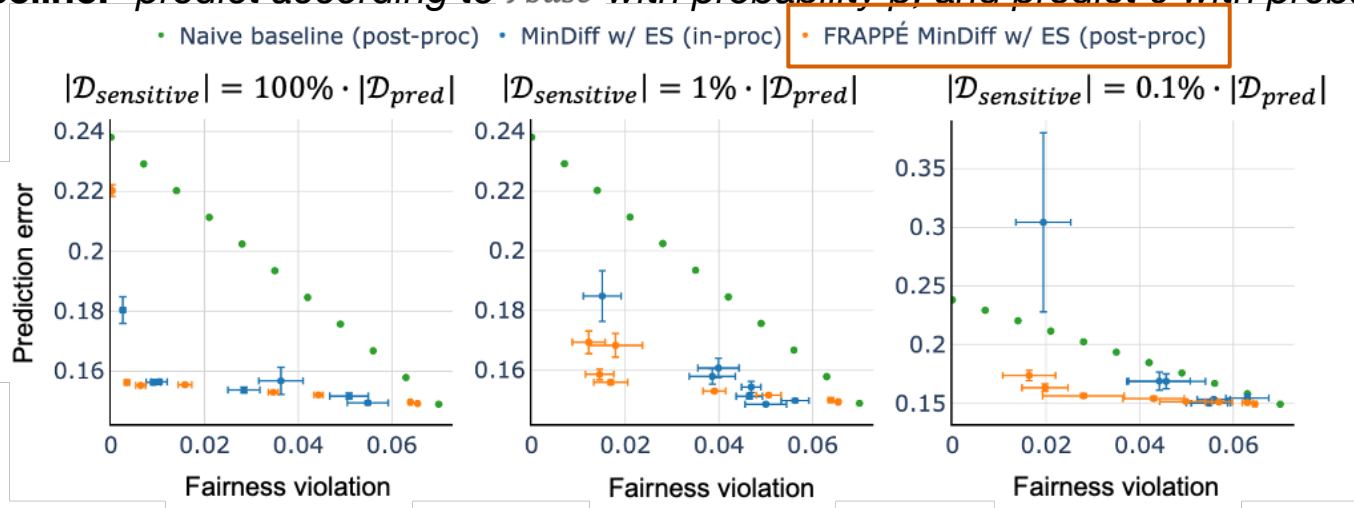
Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîfrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Setup: Equal Opportunity on Adult dataset

Naive baseline: “predict according to f_{base} with probability p , and predict 0 with probability $(1-p)$ ”



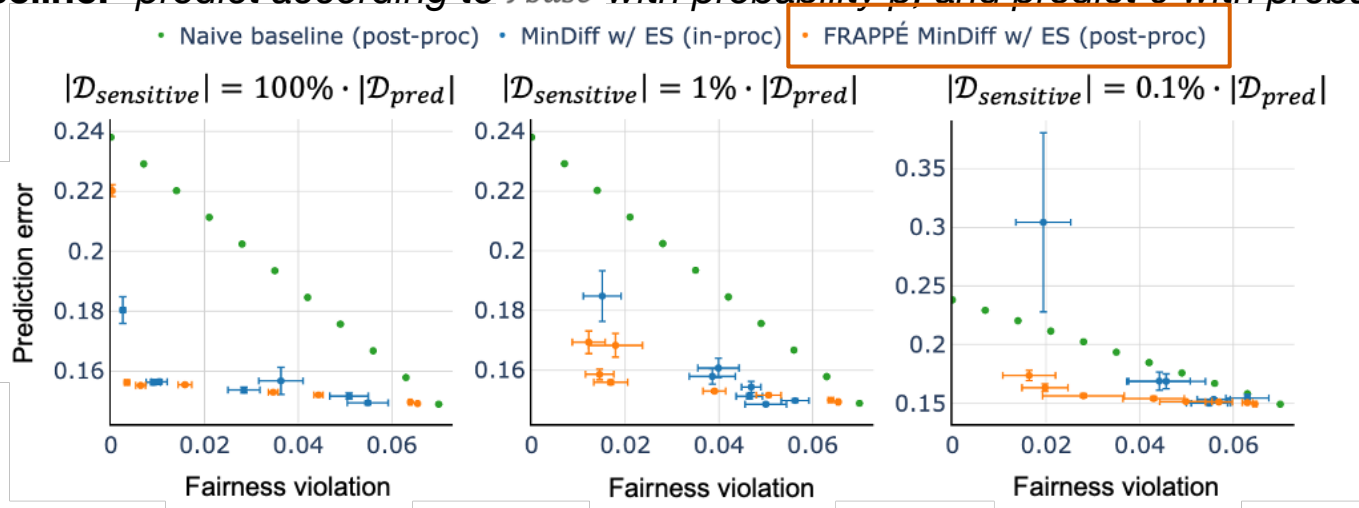
Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țifrea¹ Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Setup: Equal Opportunity on Adult dataset

Naive baseline: “predict according to f_{base} with probability p , and predict 0 with probability $(1-p)$ ”



computation time ~8x faster
than in-processing

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țifrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \underset{f}{\operatorname{argmin}} \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țifrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + T(x)$$

(logit additive for classification)

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țifrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + \underline{T(x)} \quad \text{not group-dependent}$$

(logit additive for classification)

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + \underline{T(x)} \quad \text{not group-dependent}$$

(logit additive for classification)

In-processing:

$$OPT_{IP}(f; \lambda) = \mathcal{L}_{pred}(f; \mathcal{D}_{pred}) + \lambda \mathcal{L}_{fair}(f; \mathcal{D}_{sensitive})$$

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + \underline{T(x)} \quad \text{not group-dependent}$$

(logit additive for classification)

In-processing:

$$OPT_{IP}(f; \lambda) = \mathcal{L}_{pred}(f; \mathcal{D}_{pred}) + \lambda \mathcal{L}_{fair}(f; \mathcal{D}_{sensitive})$$

Proposed post-processing for learning T :

$$OPT_{PP}(T; \lambda) = \text{Discrepancy}((f_{base} + T) \parallel f_{base}; \mathcal{D}_{unlab}) + \lambda \mathcal{L}_{fair}(f_{base} + T; \mathcal{D}_{sensitive})$$

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + \underline{T(x)} \quad \text{not group-dependent}$$

(logit additive for classification)

In-processing:

$$OPT_{IP}(f; \lambda) = \mathcal{L}_{pred}(f; \mathcal{D}_{pred}) + \lambda \mathcal{L}_{fair}(f; \mathcal{D}_{sensitive})$$

Proposed post-processing for learning T :

$$OPT_{PP}(T; \lambda) = \text{Discrepancy}((f_{base} + T) \parallel f_{base}; \mathcal{D}_{unlab}) + \underbrace{\lambda \mathcal{L}_{fair}(f_{base} + T; \mathcal{D}_{sensitive})}_{\text{any notion of fairness}}$$

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîrea^{*1} Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + \underline{T(x)} \quad \text{not group-dependent}$$

(logit additive for classification)

In-processing:

$$OPT_{IP}(f; \lambda) = \mathcal{L}_{pred}(f; \mathcal{D}_{pred}) + \lambda \mathcal{L}_{fair}(f; \mathcal{D}_{sensitive})$$

Proposed post-processing for learning T :

$$OPT_{PP}(T; \lambda) = \underbrace{Discrepancy((f_{base} + T) \parallel f_{base}; \mathcal{D}_{unlab})}_{\substack{\text{output discrepancy} \\ \text{related to } \mathcal{L}_{pred} \text{ e.g. MSE, KL divergence etc}}} + \underbrace{\lambda \mathcal{L}_{fair}(f_{base} + T; \mathcal{D}_{sensitive})}_{\text{any notion of fairness}}$$

Modular sample-efficient fairness mitigations

FRAPPÉ: A Group Fairness Framework for Post-Processing Everything

Alexandru Țîrea¹ Preethi Lahoti² Ben Packer² Yoni Halpern² Ahmad Beirami² Flavien Prost²

Accurate but unfair model:

$$f_{base} := \operatorname{argmin}_f \mathcal{L}_{pred}(f; \mathcal{D}_{pred})$$

Proposed post-hoc transformation:

$$f_{fair}(x) = f_{base}(x) + \underline{T(x)} \quad \text{not group-dependent}$$

(logit additive for classification)

In-processing:

$$OPT_{IP}(f; \lambda) = \mathcal{L}_{pred}(f; \mathcal{D}_{pred}) + \lambda \mathcal{L}_{fair}(f; \mathcal{D}_{sensitive})$$

Proposed post-processing for learning T :

$$OPT_{PP}(T; \lambda) = \underbrace{\text{Discrepancy}((f_{base} + T) \parallel f_{base}; \mathcal{D}_{unlab})}_{\text{output discrepancy related to } \mathcal{L}_{pred} \text{ e.g. MSE, KL divergence etc}} + \underbrace{\lambda \mathcal{L}_{fair}(f_{base} + T; \mathcal{D}_{sensitive})}_{\text{any notion of fairness}}$$

output discrepancy
related to \mathcal{L}_{pred} e.g. MSE, KL divergence etc

any notion of fairness

$$\mathcal{D}_{unlab} = \{x_i\}_{i=1}^N \quad \text{unlabeled data}$$

Instances of modular multi-objective learning

LLM alignment

Asymptotics of Language Model Alignment

Joy Qiping Yang
University of Sydney
Sydney, Australia
qyan6238@uni.sydney.edu.au

Salman Salamatian
Massachusetts Institute of Technology
Cambridge, MA, USA
salmansa@mit.edu

Ziteng Sun, Ananda Theertha Suresh, Ahmad Beirami
Google Research
New York, NY, USA
{zitengsun, theertha, beirami}@google.com

Out-of-domain generalization

OVERPARAMETERISATION AND WORST-CASE GENERALISATION: FRIEND OR FOE?

Aditya Krishna Menon, Ankit Singh Rawat & Sanjiv Kumar
Google Research
New York, NY
{adityakmenon, ankitsrawat, sanjivk}@google.com

Adversarial robustness

Understanding and Mitigating the Tradeoff Between Robustness and Accuracy

Aditi Raghunathan^{*1} Sang Michael Xie^{*1} Fanny Yang² John C. Duchi¹ Percy Liang¹

Unlabeled Data Improves Adversarial Robustness

Yair Carmon^{*}
Stanford University
yairc@stanford.edu

Aditi Raghunathan^{*}
Stanford University
aditir@stanford.edu

Ludwig Schmidt
UC Berkeley
ludwig@berkeley.edu

Percy Liang
Stanford University
pliang@cs.stanford.edu

John C. Duchi
Stanford University
jduchi@stanford.edu

Summary: Modular fairness mitigations

More sample efficient than in-processing

iff learning the fairness correction module is statistically efficient

e.g. $T(x)$ is not a complex function, $T(x)$ has low-dimensional structure (e.g. sparsity)

Summary: Modular fairness mitigations

More sample efficient than in-processing

iff learning the fairness correction module is statistically efficient

e.g. $T(x)$ is not a complex function, $T(x)$ has low-dimensional structure (e.g. sparsity)

Effective technique to induce any notion of fairness

iff fairness violations can be measured from observational data

e.g. $T(X)$ implicitly estimates $P(A|X)$ which might be unidentifiable from observational data

Assessing Algorithmic Fairness with Unobserved
Protected Class Using Data Combination

Nathan Kallus

Cornell University, kallus@cornell.edu

Xiaojie Mao

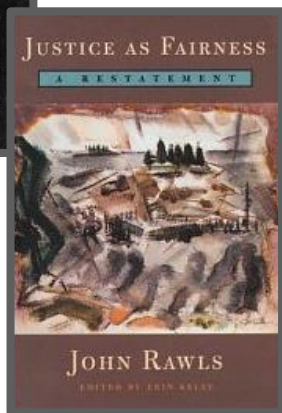
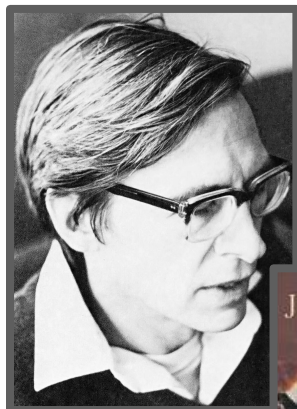
Cornell University, xm77@cornell.edu

Angela Zhou

Cornell University, az434@cornell.edu

Fairness with no group labels

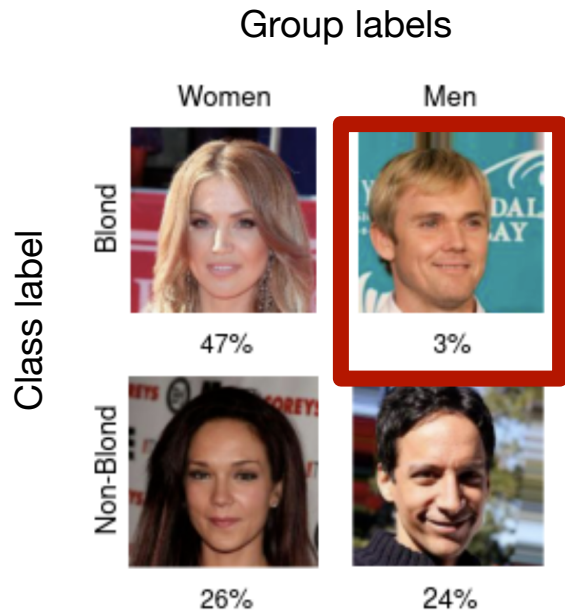
Fairness as worst-group performance



Definition A hypothesis h^* satisfies Rawlsian max-min fairness if it maximizes the accuracy of the worst-off group

$$h^* = \arg \max_h \min_{a \in \mathcal{A}} \text{Acc}(h|A = a)$$

Mitigation strategies for worst-group fairness



CelebA dataset

If we know group labels:

- importance weighting (IW)
- group distributionally robust optimization (GDRO)

DISTRIBUTIONALLY ROBUST NEURAL NETWORKS
FOR GROUP SHIFTS: ON THE IMPORTANCE OF
REGULARIZATION FOR WORST-CASE GENERALIZATION

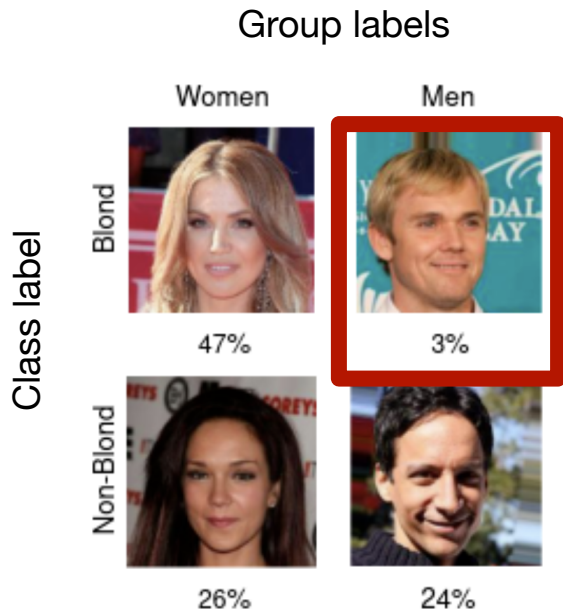
Shiori Sagawa*
Stanford University
ssagawa@cs.stanford.edu

Pang Wei Koh*
Stanford University
pangwei@cs.stanford.edu

Tatsunori B. Hashimoto
Microsoft
tahashim@microsoft.com

Percy Liang
Stanford University
pliang@cs.stanford.edu

Mitigation strategies for worst-group fairness



CelebA dataset

If we know group labels:

- importance weighting (IW)
- group distributionally robust optimization (GDRO)

DISTRIBUTIONALLY ROBUST NEURAL NETWORKS
FOR GROUP SHIFTS: ON THE IMPORTANCE OF
REGULARIZATION FOR WORST-CASE GENERALIZATION

Shiori Sagawa*
Stanford University
ssagawa@cs.stanford.edu

Pang Wei Koh*
Stanford University
pangwei@cs.stanford.edu

Tatsunori B. Hashimoto
Microsoft
tahashim@microsoft.com

Percy Liang
Stanford University
pliang@cs.stanford.edu

In the absence of group labels:

Two-stage method

- 1) identify worse-off group
- 2) employ e.g. IW/GDRO to improve worst-group error

Fairness via distributionally robust optimization (DRO)

$$\mathcal{R}_{erm}(\theta) := \mathbb{E}_P[\ell(\theta; Z)]$$

Fairness Without Demographics in Repeated Loss Minimization

Tatsunori B. Hashimoto^{1,2} Megha Srivastava¹ Hongseok Namkoong³ Percy Liang¹

Fairness via distributionally robust optimization (DRO)

$$\mathcal{R}_{erm}(\theta) := \mathbb{E}_P[\ell(\theta; Z)]$$

$$\mathcal{R}_{dro}(\theta; r) := \sup_{Q \in \mathcal{B}(P, r)} \mathbb{E}_Q[\ell(\theta; Z)] \text{ *worst-case loss wrt the uncertainty set } Q*$$

Fairness Without Demographics in Repeated Loss Minimization

Tatsunori B. Hashimoto^{1,2} Megha Srivastava¹ Hongseok Namkoong³ Percy Liang¹

Fairness via distributionally robust optimization (DRO)

Fairness Without Demographics in Repeated Loss Minimization
Tatsunori B. Hashimoto ^{1,2} Megha Srivastava ¹ Hongseok Namkoong ³ Percy Liang ¹

$$\mathcal{R}_{erm}(\theta) := \mathbb{E}_P[\ell(\theta; Z)]$$

$$\mathcal{R}_{dro}(\theta; r) := \sup_{Q \in \mathcal{B}(P, r)} \mathbb{E}_Q[\ell(\theta; Z)] \text{ *worst-case loss wrt the uncertainty set } Q*$$

P = (marginal)
data distribution



Fairness via distributionally robust optimization (DRO)

Fairness Without Demographics in Repeated Loss Minimization

Tatsunori B. Hashimoto^{1,2} Megha Srivastava¹ Hongseok Namkoong³ Percy Liang¹

$$\mathcal{R}_{erm}(\theta) := \mathbb{E}_P[\ell(\theta; Z)]$$

$$\mathcal{R}_{dro}(\theta; r) := \sup_{Q \in \mathcal{B}(P, r)} \mathbb{E}_Q[\ell(\theta; Z)] \quad \text{worst-case loss wrt the uncertainty set } \mathcal{Q}$$

P = (marginal)
data distribution

r = radius of
uncertainty set

Fairness via distributionally robust optimization (DRO)

Fairness Without Demographics in Repeated Loss Minimization

Tatsunori B. Hashimoto^{1,2} Megha Srivastava¹ Hongseok Namkoong³ Percy Liang¹

$$\mathcal{R}_{erm}(\theta) := \mathbb{E}_P[\ell(\theta; Z)]$$

$$\mathcal{R}_{dro}(\theta; r) := \sup_{Q \in \mathcal{B}(P, r)} \mathbb{E}_Q[\ell(\theta; Z)] \quad \text{worst-case loss wrt the uncertainty set } Q$$

P = (marginal)
data distribution

r = radius of
uncertainty set

determined by

α_{\min} minority
group proportion

Fairness via distributionally robust optimization (DRO)

Fairness Without Demographics in Repeated Loss Minimization

Tatsunori B. Hashimoto^{1,2} Megha Srivastava¹ Hongseok Namkoong³ Percy Liang¹

$$\mathcal{R}_{erm}(\theta) := \mathbb{E}_P[\ell(\theta; Z)]$$

$$\mathcal{R}_{dro}(\theta; r) := \sup_{Q \in \mathcal{B}(P, r)} \mathbb{E}_Q[\ell(\theta; Z)] \quad \text{worst-case loss wrt the uncertainty set } \mathcal{Q}$$

P = (marginal)
data distribution

r = radius of
uncertainty set

determined by

α_{\min} minority
group proportion

What if no group labels available?

A: pick a lower bound for α_{\min}

Detect worst-group using a biased classifier

DRO: upweights high-loss samples.

Alternative: Two-stage method

- 1) use **biased** classifier to identify *error set*
- 2) train fair classifier via IW / GroupDRO

Detect worst-group using a biased classifier

DRO: upweights high-loss samples.

Alternative: Two-stage method

- 1) use **biased** classifier to identify *error set*
- 2) train fair classifier via IW / GroupDRO

Why are two-stage methods expected to work?



Majority group

Intuition: a biased classifier will predict based on the stronger correlation.
e.g. background

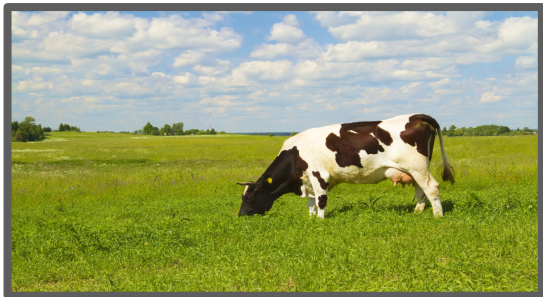
Detect worst-group using a biased classifier

DRO: upweights high-loss samples.

Alternative: Two-stage method

- 1) use **biased** classifier to identify *error set*
- 2) train fair classifier via IW / GroupDRO

Why are two-stage methods expected to work?



Majority group



Minority group

Intuition: a biased classifier will predict based on the stronger correlation.
e.g. background



incorrect predictions where
spurious correlation does not hold
i.e. minority groups

How to train a biased classifier?

Setting 1: group labels available for validation set

How to train a biased classifier?

Setting 1: group labels available for validation set

Examples:

- heavy regularization
e.g. via early stopping
- custom loss function
e.g. amplify “easy” examples

Just Train Twice: Improving Group Robustness without Training Group Information

Evan Zheran Liu^{*1} Behzad Haghgoo^{*1} Annie S. Chen^{*1} Aditi Raghunathan¹ Pang Wei Koh¹
Shiori Sagawa¹ Percy Liang¹ Chelsea Finn¹

Learning from Failure: Training Debiased Classifier from Biased Classifier

Junhyun Nam¹ Hyuntak Cha² Sungsoo Ahn¹ Jaeho Lee¹ Jinwoo Shin^{1,2}
¹School of Electrical Engineering, KAIST
²Graduate School of AI, KAIST
{junhyun.nam, hyuntak.cha, sungsoo.ahn, jaeho-lee, jinwoos}@kaist.ac.kr

Use worst-group validation error to select regularization strength, IW weights etc.

How to train a biased classifier?

Setting 2: no group labels at all

How to train a biased classifier?

Setting 2: no group labels at all

Examples:

- identify groups from training AND validation data with ensemble of biased classifiers to reduce noise
- post-hoc logit adjustment using $P\left(Y|\hat{Y}_{biased}\right)$ as an estimate of $P(Y|A)$

Boosting worst-group accuracy without any group annotations

Vincent Bardenhagen*, Alexandru Tifrea*, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland
{vbardenha,tifreaa,fan.yang}@ethz.ch

Group Robust Classification Without Any Group Information

Christos Tsirigotis* Joao Monteiro Pau Rodriguez
Université de Montréal, Mila, ServiceNow Research ServiceNow Research Apple MLR

David Vazquez Aaron Courville†
ServiceNow Research Université de Montréal, Mila, CIFAR CAI Chair

How to train a biased classifier?

Setting 2: no group labels at all

Examples:

- identify groups from training AND validation data with ensemble of biased classifiers to reduce noise
- post-hoc logit adjustment using $P(Y|\hat{Y}_{biased})$ as an estimate of $P(Y|A)$

**Boosting worst-group accuracy
without any group annotations**

Vincent Bardenhagen*, Alexandru Tifrea*, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland
{vbardenha,tifreaa,fan.yang}@ethz.ch

**Group Robust Classification
Without Any Group Information**

Christos Tsirigotis* Joao Monteiro Pau Rodriguez
Université de Montréal, Mila, ServiceNow Research ServiceNow Research Apple MLR

David Vazquez Aaron Courville†
ServiceNow Research Université de Montréal, Mila, CIFAR CAI Chair

	Tuning	Corrupt-MNIST		Waterbirds		CelebA		Color MNIST		Adult		Poverty	
		Avg	Wg	Avg	Wg	Avg	Wg	Avg	Wg	Avg	Wg	Avg	Wg
No group labels	ERM	99.6	71.2	97.9	74.9	94.3	60.7	99.8	82.6	80.1	41.6	87.6	55.6
	Ours	99.0	96.5	97.5	78.5	88.0	78.9	99.3	96.6	81.2	68.0	86.3	50.0
Val group labels	ERM WG	99.5	79.8	97.6	86.7	93.1	77.8	99.7	84.4	78.9	61.2	87.7	51.5
	JTT	99.1	91.3	93.3	86.7	88.0	81.1	98.3	94.8	77.8	63.3	64.5	60.5

Similar average and worst-group accuracy for two-stage methods:

- with no group labels
- with validation group labels

DRO mitigations in the presence of outliers

Recall: Two-stage methods

- 1) use biased classifier to identify *error set*
- 2) train fair classifier via IW / GDRO

DRO mitigations in the presence of outliers

Recall: Two-stage methods

- 1) use biased classifier to identify *error set*
- 2) train fair classifier via IW / GDRO



The opposite of what robust statistics literature recommends!
e.g. can amplify outliers, noisy samples etc

Can we get both fairness and robustness to outliers?

DRO mitigations in the presence of outliers

Recall: Two-stage methods

- 1) use biased classifier to identify *error set*
- 2) train fair classifier via IW / GDRO



The opposite of what robust statistics literature recommends!
e.g. can amplify outliers, noisy samples etc

Can we get both fairness and robustness to outliers?

Robust Mixture Learning when Outliers Overwhelm Small Groups

Daniil Dmitriev¹, Rares-Darius Buhai¹, Stefan Tiegel¹, Alexander Wolters², Gleb Novikov³, Amartya Sanyal⁴, David Steurer¹, and Fanny Yang¹

Clustering algorithm that is

- applicable even for $|\text{Outliers}| \gg |\text{Minority group}|$
- computationally efficient
- information-theoretically optimal

DRO mitigations in the presence of outliers

Recall: Two-stage methods

- 1) use biased classifier to identify *error set*
- 2) train fair classifier via IW / GDRO



The opposite of what robust statistics literature recommends!
e.g. can amplify outliers, noisy samples etc

Can we get both fairness and robustness to outliers?

Robust Mixture Learning when Outliers Overwhelm Small Groups

Daniil Dmitriev¹, Rares-Darius Buhai¹, Stefan Tiegel¹, Alexander Wolters², Gleb Novikov³, Amartya Sanyal⁴, David Steurer¹, and Fanny Yang¹

Clustering algorithm that is

- applicable even for $|\text{Outliers}| \gg |\text{Minority group}|$
- computationally efficient
- information-theoretically optimal

Fairness without Demographics through Adversarially Reweighted Learning

Preethi Lahoti *
plahoti@mpi-inf.mpg.de
Max Planck Institute for Informatics

Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost,
Nithum Thain, Xuezhong Wang, Ed H. Chi
Google Research

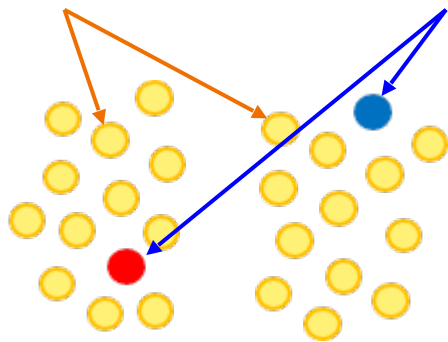
Idea: only upweight samples in the error set that are computationally identifiable using simple function class \mathcal{F} .

Fairness in the low-label regime

Low-label regime

unlabeled data

labeled data

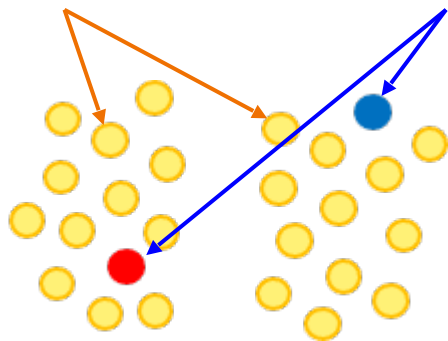


Low-label regime

Research questions

unlabeled data

labeled data



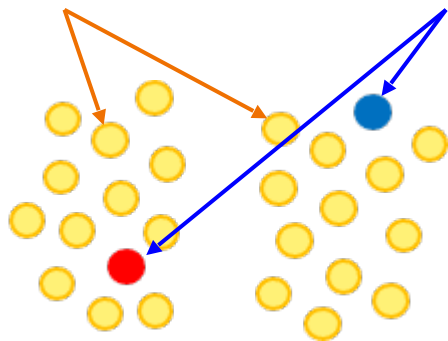
- 1) How to acquire the labeled data?
- 2) How to learn from both labeled and unlabeled data?

Low-label regime

Research questions

unlabeled data

labeled data



- 1) How to acquire the labeled data?
- 2) How to learn from both labeled and unlabeled data?

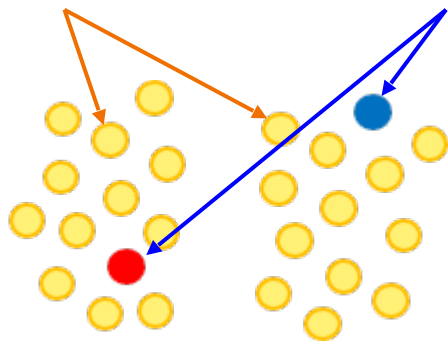
active learning

Low-label regime

Research questions

unlabeled data

labeled data



- 1) How to acquire the labeled data?
- 2) How to learn from both labeled and unlabeled data?

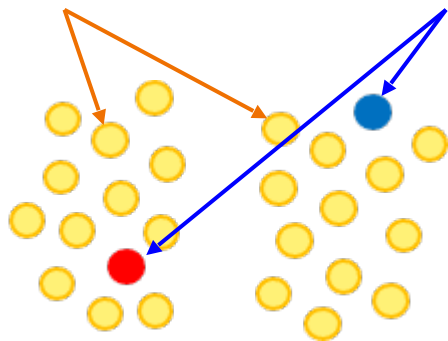
active learning

semi-supervised learning

Low-label regime

unlabeled data

labeled data



Research questions

- 1) How to acquire the labeled data?
- 2) How to learn from both labeled and unlabeled data?

active learning

semi-supervised learning

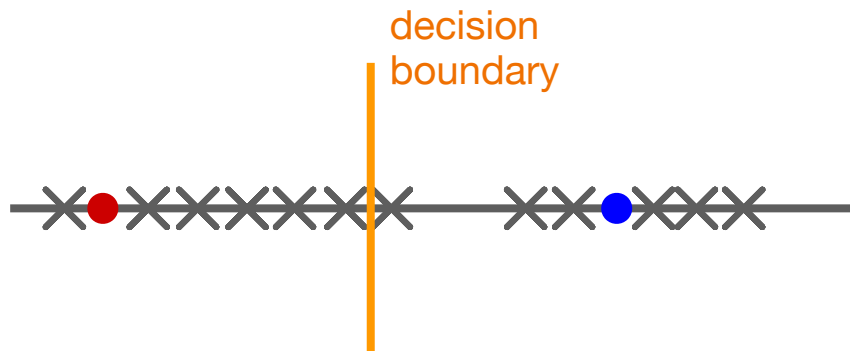
Fairness problems

- class imbalance
- group imbalance
(but potentially balanced classes)

Uncertainty sampling-based active learning

Uncertainty sampling

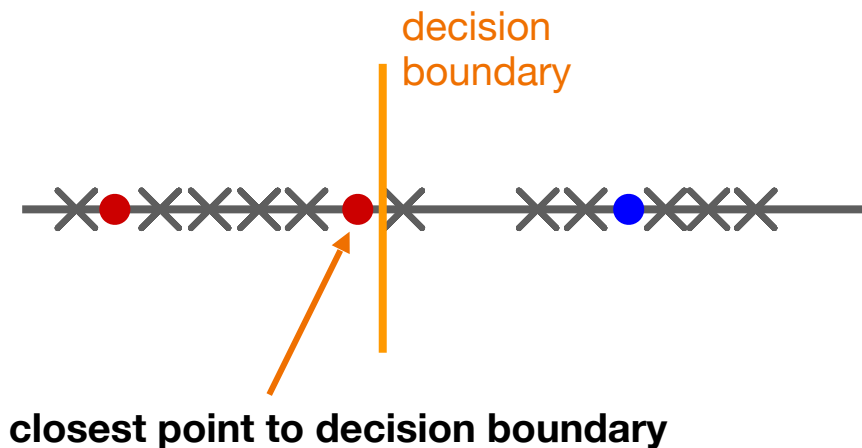
“binary search to find decision boundary”



Uncertainty sampling-based active learning

Uncertainty sampling

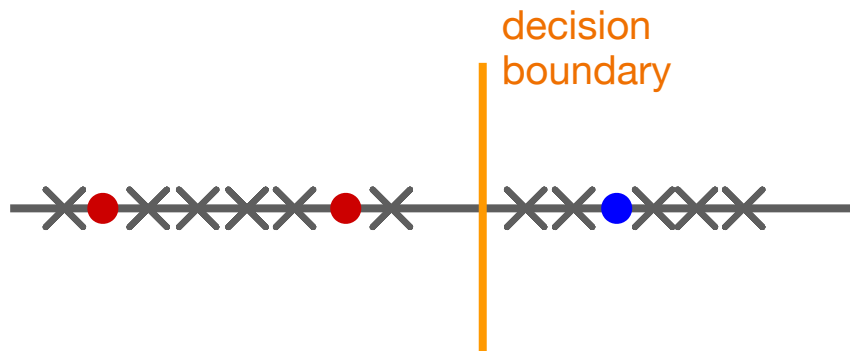
“binary search to find decision boundary”



Uncertainty sampling-based active learning

Uncertainty sampling

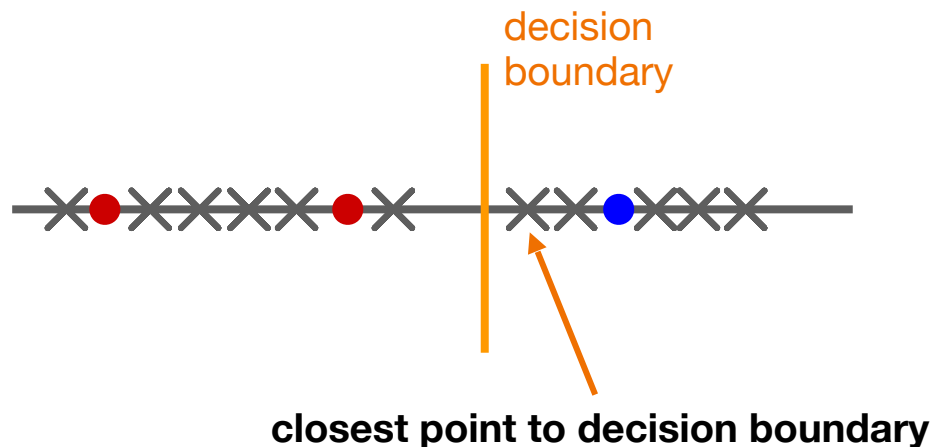
“binary search to find decision boundary”



Uncertainty sampling-based active learning

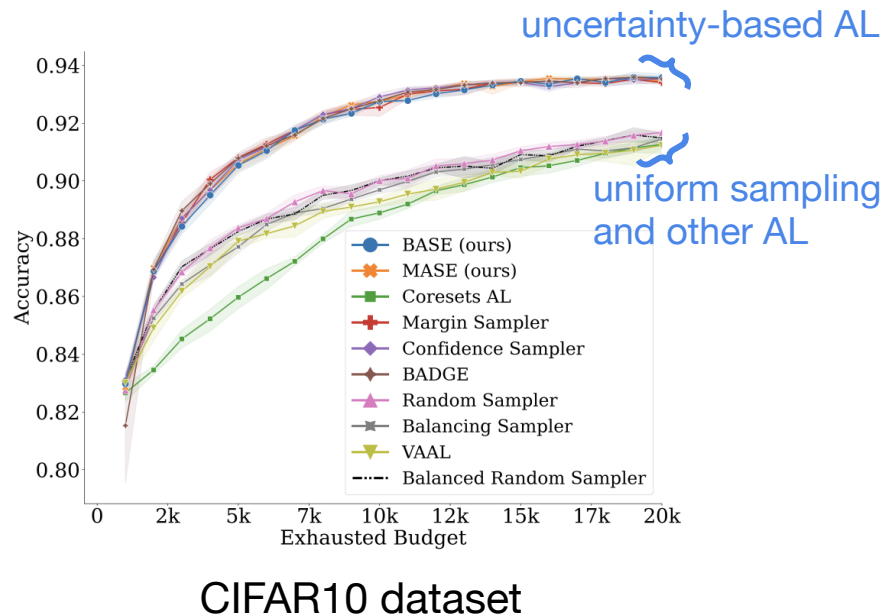
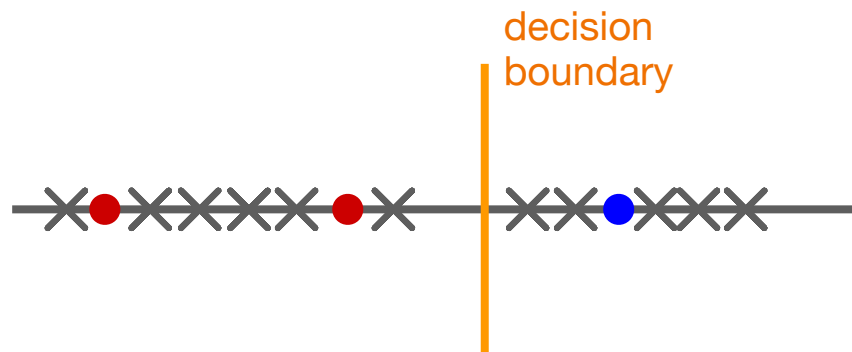
Uncertainty sampling

“binary search to find decision boundary”



Uncertainty sampling-based active learning

Uncertainty sampling
“binary search to find decision boundary”



U-AL is more label efficient than uniform sampling or other AL

Standard active learning can improve fairness

Class-imbalanced classification

Learning on the Border:
Active Learning in Imbalanced Data Classification

Şeyda Ertekin¹, Jian Huang², Léon Bottou³, C. Lee Giles^{2,1}



Focus on [linear classification](#)

Standard active learning can improve fairness

Class-imbalanced classification

true decision boundary best **avg-case** classifier



Focus on **linear classification**

Learning on the Border:
Active Learning in Imbalanced Data Classification

Şeyda Ertekin¹, Jian Huang², Léon Bottou³, C. Lee Giles^{2,1}

Standard active learning can improve fairness

Class-imbalanced classification

true decision boundary

best **avg-case** classifier



Focus on **linear classification**

Learning on the Border:
Active Learning in Imbalanced Data Classification

Şeyda Ertekin¹, Jian Huang², Léon Bottou³, C. Lee Giles^{2,1}



Decision boundary of biased classifier is closer to minority class

Standard active learning can improve fairness

Class-imbalanced classification

Learning on the Border:
Active Learning in Imbalanced Data Classification

Şeyda Ertekin¹, Jian Huang², Léon Bottou³, C. Lee Giles^{2,1}

true decision boundary best **avg-case** classifier



Decision boundary of biased classifier is closer to minority class



U-AL tends to select more minority points to be labeled

Focus on **linear classification**

Standard active learning can improve fairness

Class-imbalanced classification

true decision boundary best **avg-case** classifier



Focus on **linear classification**

Learning on the Border:
Active Learning in Imbalanced Data Classification

Şeyda Ertekin¹, Jian Huang², Léon Bottou³, C. Lee Giles^{2,1}



Decision boundary of biased classifier is closer to minority class



U-AL tends to select more minority points to be labeled



U-AL collects a more balanced labeled set

Standard active learning can improve fairness

Class-imbalanced classification

U-AL also mitigates class imbalance in non-linear classification!

Active Learning at the ImageNet Scale

Zeyad Ali Sami Emam^{*†‡} Hong-Min Chu[†] Ping-Yeh Chiang[†] Wojciech Czaja[†]
 zeyad@umd.edu hmchu@umd.edu pchiang@umd.edu wojtek@umd.edu

Richard Leapman[‡] Micah Goldblum[§] Tom Goldstein[†]
 leapmanr@mail.nih.gov goldblum@nyu.edu tomg@umd.edu

Improving class and group imbalanced classification with uncertainty-based active learning

Alexandru Tifrea^{*} TIFREAA@INF.ETHZ.CH
 Department of Computer Science, ETH Zurich

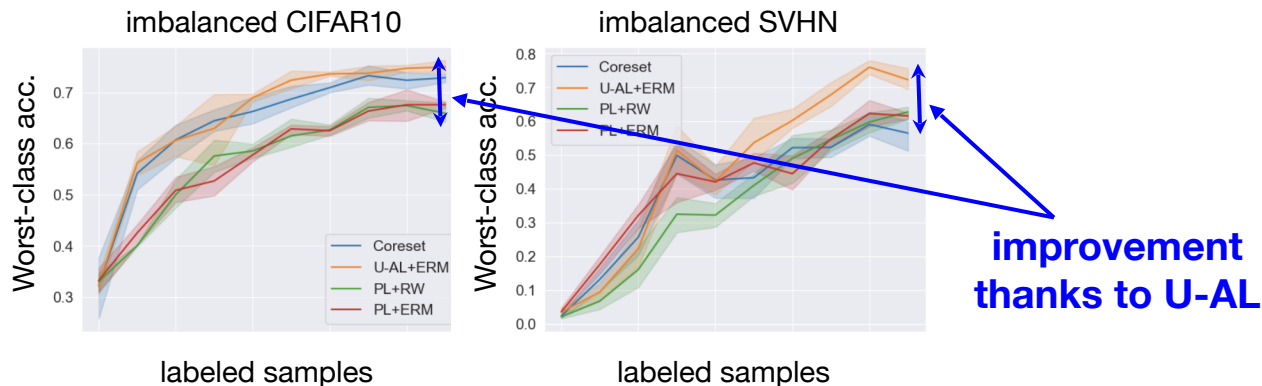
John Hill^{*} JHILL326@GATECH.EDU
 Department of Computer Science, Georgia Institute of Technology

Fanny Yang FAN.YANG@INF.ETHZ.CH
 Department of Computer Science, ETH Zurich

Algorithm Selection for Deep Active Learning with Imbalanced Datasets

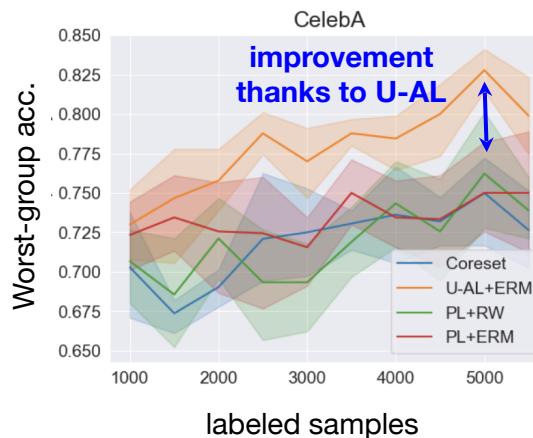
Jifan Zhang Shuai Shao Saurabh Verma
 University of Wisconsin - Madison Meta Inc. Meta Inc.
 Madison, WI 53715 Menlo Park, CA 94025 Menlo Park, CA 94025
 jifan@cs.wisc.edu sshao@meta.com saurabh08@meta.com

Robert Nowak
 University of Wisconsin - Madison
 Madison, WI 53715
 rdnowak@wisc.edu



Standard active learning can improve fairness

Group-imbalanced classification



Improving class and group imbalanced classification with uncertainty-based active learning

Alexandru Tifrea*

Department of Computer Science, ETH Zurich

TIFREAA@INF.ETHZ.CH

John Hill*

Department of Computer Science, Georgia Institute of Technology

JHILL326@GATECH.EDU

Fanny Yang

Department of Computer Science, ETH Zurich

FAN.YANG@INF.ETHZ.CH

CAN ACTIVE LEARNING PREEMPTIVELY MITIGATE FAIRNESS ISSUES?

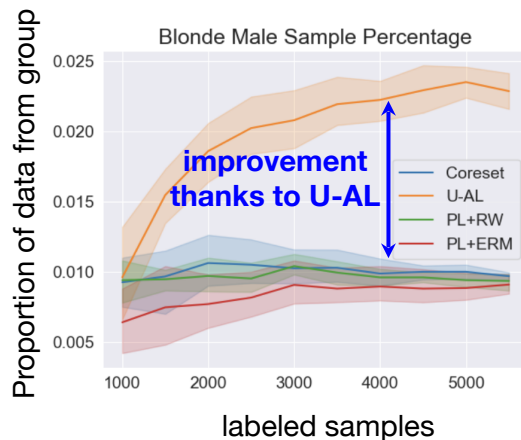
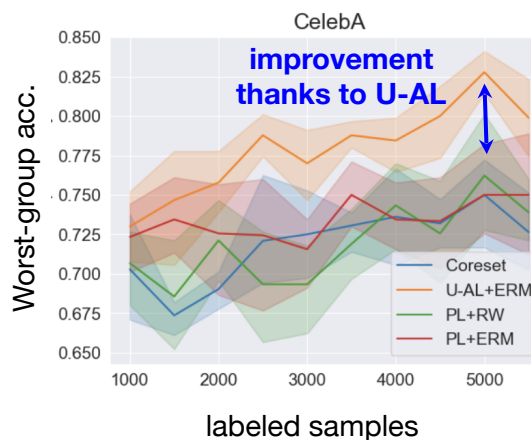
Frédéric Branchaud-Charron*, Parmida Atighehchian*, Pau Rodríguez, Grace Abuhamad, Alexandre Lacoste

ServiceNow

{fr.branchaud-charron, parmida.atighehchian}@servicenow.com

Standard active learning can improve fairness

Group-imbalanced classification



Improving class and group imbalanced classification
with uncertainty-based active learning

Alexandru Tifrea*

Department of Computer Science, ETH Zurich

TIFREAA@INF.ETHZ.CH

John Hill*

Department of Computer Science, Georgia Institute of Technology

JHILL326@GATECH.EDU

Fanny Yang

Department of Computer Science, ETH Zurich

FAN.YANG@INF.ETHZ.CH

CAN ACTIVE LEARNING PREEMPTIVELY MITIGATE
FAIRNESS ISSUES?

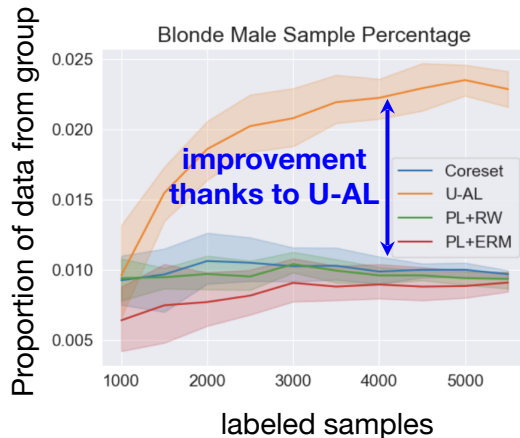
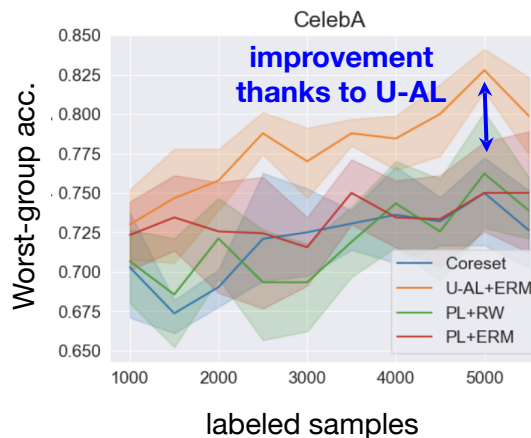
**Frédéric Branchaud-Charron*, Parmida Atighehchian*, Pau Rodríguez,
Grace Abuhamad, Alexandre Lacoste**

ServiceNow

{fr.branchaud-charron, parmida.atighehchian}@servicenow.com

Standard active learning can improve fairness

Group-imbalanced classification



Improving class and group imbalanced classification
with uncertainty-based active learning

Alexandru Tifrea*

Department of Computer Science, ETH Zurich

TIFREAA@INF.ETHZ.CH

John Hill*

Department of Computer Science, Georgia Institute of Technology

JHILL326@GATECH.EDU

Fanny Yang

Department of Computer Science, ETH Zurich

FAN.YANG@INF.ETHZ.CH

CAN ACTIVE LEARNING PREEMPTIVELY MITIGATE
FAIRNESS ISSUES?

**Frédéric Branchaud-Charron*, Parmida Atighehchian*, Pau Rodríguez,
Grace Abuhamad, Alexandre Lacoste**

ServiceNow

{fr.branchaud-charron, parmida.atighehchian}@servicenow.com

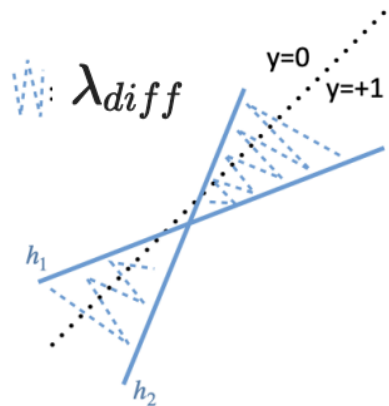
Takeaways

- no explicit group information used anywhere during sampling/learning!
- not *all* AL strategies help (e.g. coreset sampling)
- U-AL+ERM can be better than passive learning + reweighting

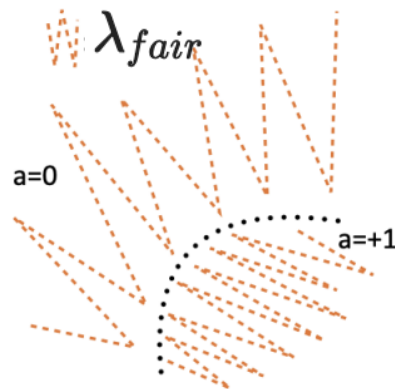
Using group labels for active learning

Acquire labels for samples selected according to:

$$P_{AL}(X) \sim \frac{1}{2} \lambda_{diff}(X) + \frac{1}{2} \lambda_{fair}(X)$$



Informativeness criterion:
Disagreement region of ensemble

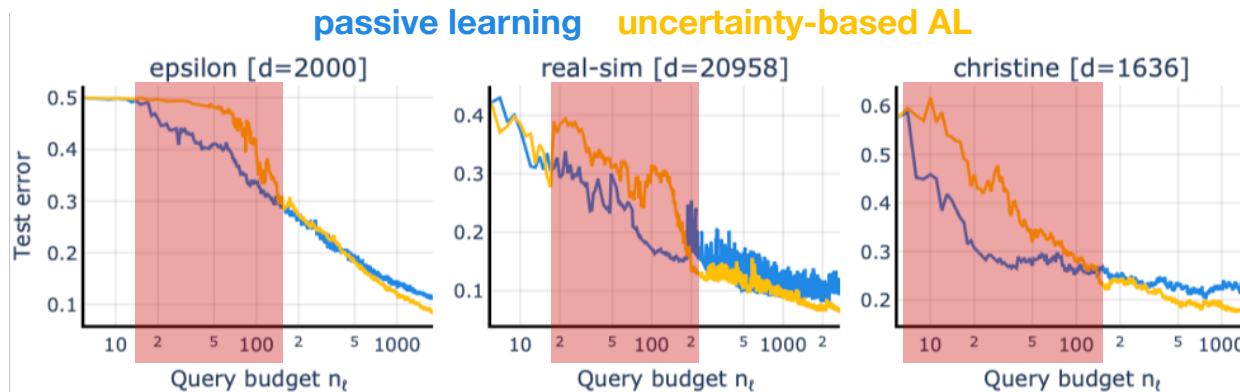


Fairness criterion:
Uniform mass on all groups

Fair Active Learning in Low-Data Regimes

Romain Camilleri, Andrew Wagenmaker, Jamie Morgenstern, Lalit Jain, Kevin Jamieson
University of Washington, Seattle, WA
{camilr, ajwagen, jamiemmt, jamieson}@cs.washington.edu, lalitj@uw.edu

Limitations of uncertainty-based AL



Err[U-AL] > Err[PL]

U-AL can be on par with or even worse than passive learning

- For high-dimensional data
- For data with lots of label noise

**Margin-based sampling in high dimensions:
When being active is less efficient than staying passive**

Alexandru Tifrea^{*1} Jacob Clarysse^{*1} Fanny Yang¹

**On the Relationship between Data Efficiency and Error
for Uncertainty Sampling**

Stephen Mussmann¹ Percy Liang¹

Summary

A few examples of fair learning algorithms that

- (1) Have fewer data requirements than standard fairness mitigations
- (2) Leverage unlabeled data to improve fairness

Open questions

- Impact of class/group label noise
- Interplay between fairness and other evaluation metrics, beyond accuracy

 coming up in the next part

Privacy in Machine Learning

Privacy in Machine Learning

Privacy in Machine Learning



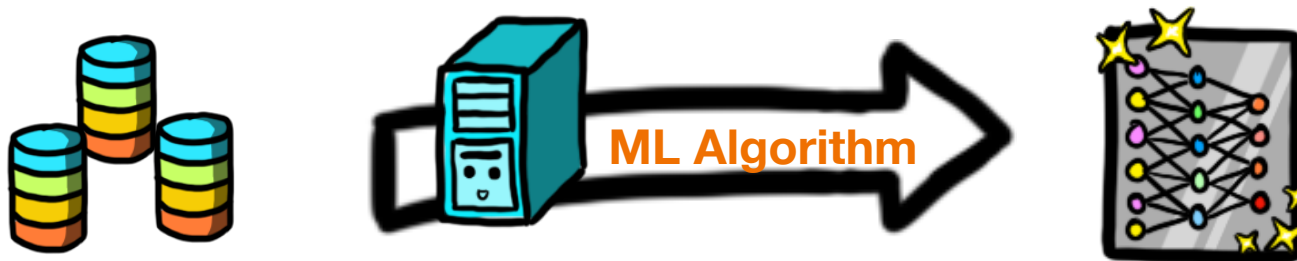
Privacy in Machine Learning



Privacy in Machine Learning

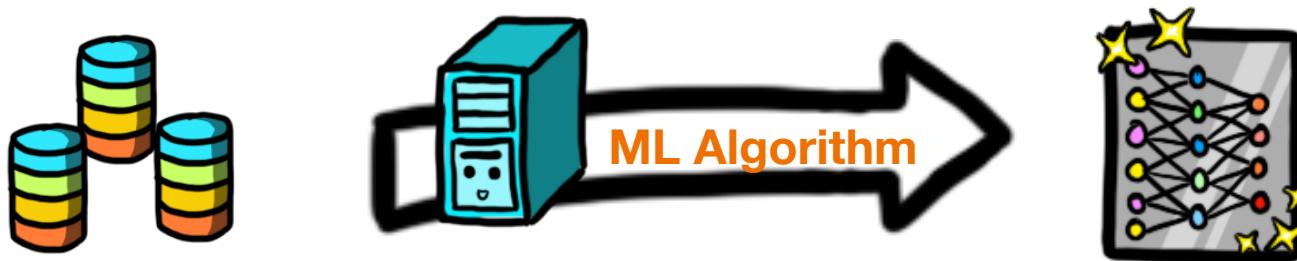


Privacy in Machine Learning



Privacy can mean a lot of things but two things are important to define:

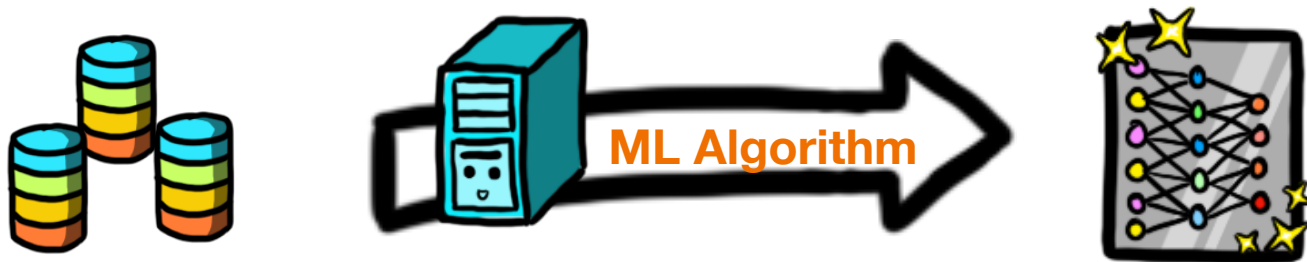
Privacy in Machine Learning



Privacy can mean a lot of things but two things are important to define:

- What is the private entity ?

Privacy in Machine Learning



Privacy can mean a lot of things but two things are important to define:

- What is the private entity ?
- What can the privacy adversary observe ?

PETs in Machine Learning



PETs in Machine Learning



- **Differential Privacy** prevents leakage of *training data from the trained model*

PETs in Machine Learning



- **Differential Privacy** prevents leakage of *training data from the trained model*
- **Multi-Party Computation** allows multiple data holders to collaboratively execute a computation without learn too much about other parties' data

PETs in Machine Learning



- **Differential Privacy** prevents leakage of *training data from the trained model*
- **Multi-Party Computation** allows multiple data holders to collaboratively execute a computation without learn too much about other parties' data
- **Fully Homomorphic encryption** based methods allow training or testing on encrypted data without decrypting the data.

PETs in Machine Learning

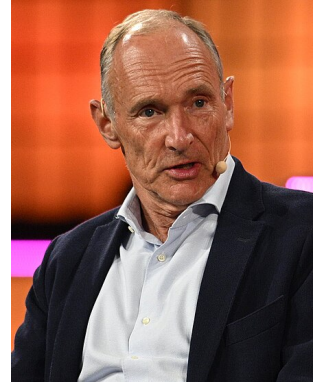


- **Differential Privacy** prevents leakage of *training data from the trained model*
- **Multi-Party Computation** allows multiple data holders to collaboratively execute a computation without learn too much about other parties' data
- **Fully Homomorphic encryption** based methods allow training or testing on encrypted data without decrypting the data.
- Many more PETs like Trusted Execution Environments (TEE), Contextual Integrity etc.

PETs in Machine Learning



- **Differential Privacy** prevents leakage of *training data from the trained model*
- **Multi-Party Computation** allows multiple data holders to collaboratively execute a computation without learn too much about other parties' data
- **Fully Homomorphic encryption** based methods allow training or testing on encrypted data without decrypting the data.
- Many more PETs like Trusted Execution Environments (TEE), Contextual Integrity etc.



“Data is a precious thing and will last longer than the systems themselves”

Sir Tim-Berners Lee

US Census and Privacy

Vulnerability of sparse data

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH
CERTAINTY?

Aloni Cohen and JN Matthews

University of Chicago

US Census and Privacy

Vulnerability of sparse data

- 2010 US Census privacy protections were vulnerable to reconstruction attacks.

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH
CERTAINTY?

Aloni Cohen and JN Matthews

University of Chicago

US Census and Privacy

Vulnerability of sparse data

- 2010 US Census privacy protections were vulnerable to reconstruction attacks.
- Some groups are disproportionately affected: 80% of NHPI (Native Hawaiian & Pacific Islander) responses in NC were fully reconstructed.

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH CERTAINTY?

Aloni Cohen and JN Matthews
University of Chicago

US Census and Privacy

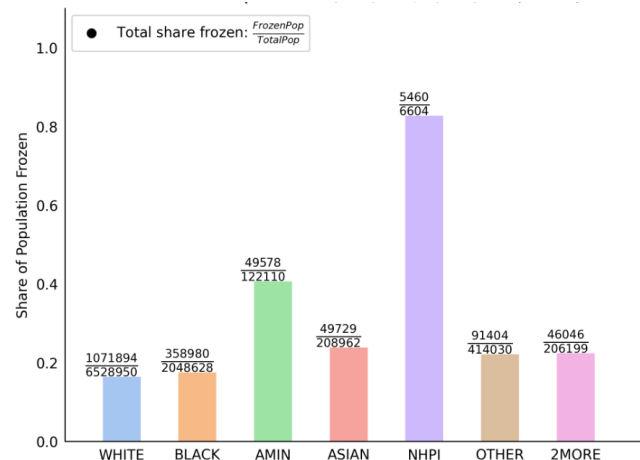
Vulnerability of sparse data

- 2010 US Census privacy protections were vulnerable to reconstruction attacks.
- Some groups are disproportionately affected: 80% of NHPI (Native Hawaiian & Pacific Islander) responses in NC were fully reconstructed.

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH CERTAINTY?

Aloni Cohen and JN Matthews

University of Chicago



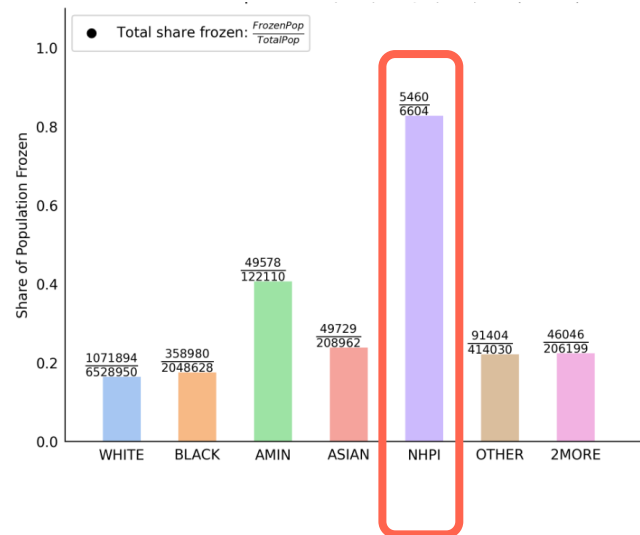
US Census and Privacy

Vulnerability of sparse data

- 2010 US Census privacy protections were vulnerable to reconstruction attacks.
- Some groups are disproportionately affected: 80% of NHPI (Native Hawaiian & Pacific Islander) responses in NC were fully reconstructed.

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH CERTAINTY?

Aloni Cohen and JN Matthews
University of Chicago



US Census and Privacy

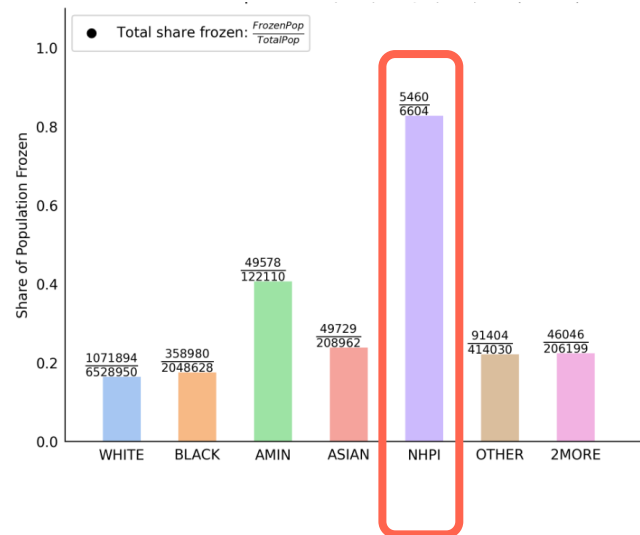
Vulnerability of sparse data

- 2010 US Census privacy protections were vulnerable to reconstruction attacks.
- Some groups are disproportionately affected: **80% of NHPI (Native Hawaiian & Pacific Islander)** responses in NC were fully reconstructed.
- Similar inferences were also shown about age and more accurate in smaller “blocks”

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH CERTAINTY?

Aloni Cohen and JN Matthews

University of Chicago



US Census and Privacy

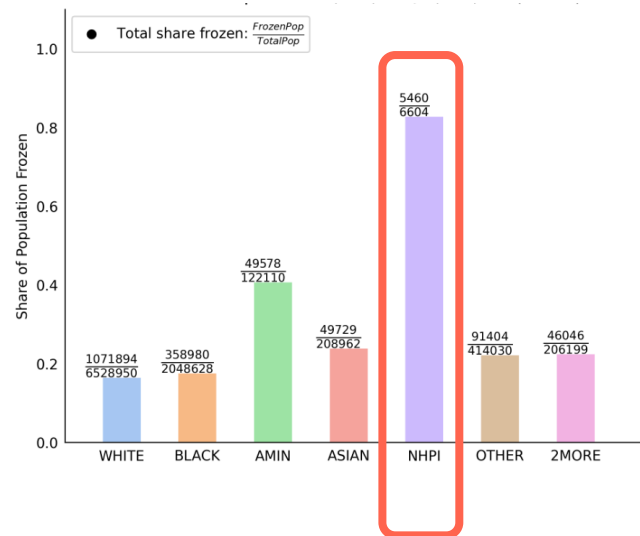
Vulnerability of sparse data

- 2010 US Census privacy protections were vulnerable to reconstruction attacks.
- Some groups are disproportionately affected: **80% of NHPI (Native Hawaiian & Pacific Islander)** responses in NC were fully reconstructed.
- Similar inferences were also shown about age and more accurate in smaller “blocks”

WHOSE 2010 CENSUS RESPONSES CAN BE RECONSTRUCTED WITH CERTAINTY?

Aloni Cohen and JN Matthews

University of Chicago



Takeaway: Often privacy violations are stronger in smaller communities.

Cost of Privacy

Cost of Privacy

Informal Theorem: If you try to answer too many questions too accurately about a dataset, there's a clever way for an attacker to piece together (almost) the entire original data.

Cost of Privacy

Informal Theorem: If you try to answer too many questions too accurately about a dataset, there's a clever way for an attacker to piece together (almost) the entire original data.

If the original dataset's privacy is to be protected, some accuracy needs to be sacrificed. The study of DP tries to control this trade-off.

Making an Algorithm Differentially Private

Making an Algorithm Differentially Private

Differential Privacy

Making an Algorithm Differentially Private

Differential Privacy

- **Differential Privacy** noises the algorithm's output to limit the exposure of any single data point

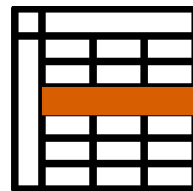
Making an Algorithm Differentially Private

Differential Privacy

- **Differential Privacy** noises the algorithm's output to limit the exposure of any single data point
- A **Differentially Private** ML algorithm produces similar models irrespective of whether Alice's data is in the dataset or Bob's

Making an Algorithm Differentially Private

Differential Privacy



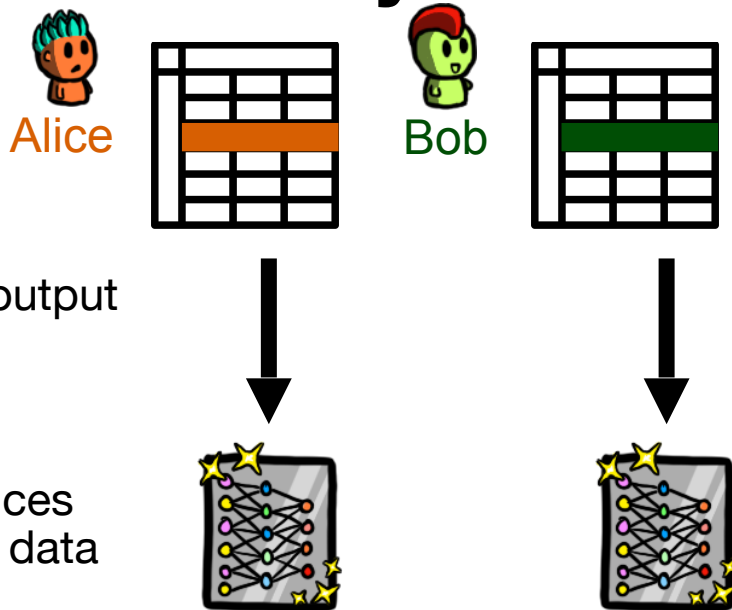
- **Differential Privacy** noises the algorithm's output to limit the exposure of any single data point
- A **Differentially Private** ML algorithm produces similar models irrespective of whether Alice's data is in the dataset or Bob's



Making an Algorithm Differentially Private

Differential Privacy

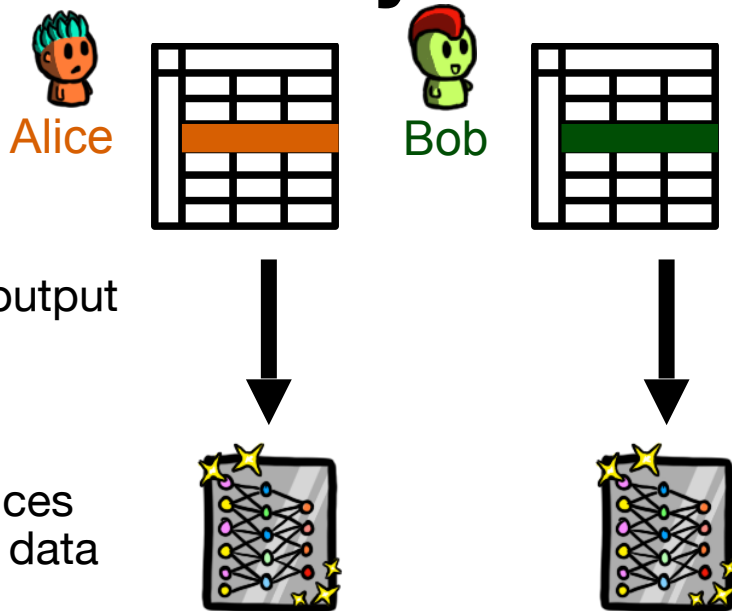
- **Differential Privacy** noises the algorithm's output to limit the exposure of any single data point
- A **Differentially Private** ML algorithm produces similar models irrespective of whether Alice's data is in the dataset or Bob's



Making an Algorithm Differentially Private

Differential Privacy

- **Differential Privacy** noises the algorithm's output to limit the exposure of any single data point
- A **Differentially Private** ML algorithm produces similar models irrespective of whether Alice's data is in the dataset or Bob's



The replacement of a single data record minimally impacts the trained model

How to make Machine Learning Private

How to make Machine Learning Private

Differential Privacy (Defn.)

How to make Machine Learning Private

Differential Privacy (Defn.)

Consider any

- Neighbouring datasets S_1 and S_2
- Output set Q

Then Algorithm is (ϵ, δ) -DP if

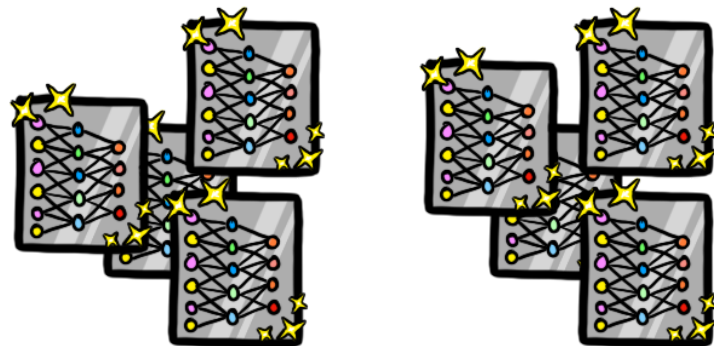
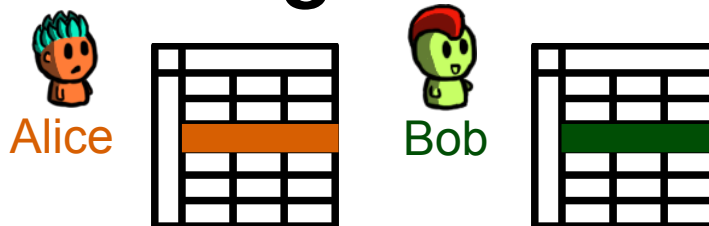
How to make Machine Learning Private

Differential Privacy (Defn.)

Consider any

- Neighbouring datasets S_1 and S_2
- Output set Q

Then Algorithm is (ϵ, δ) -DP if



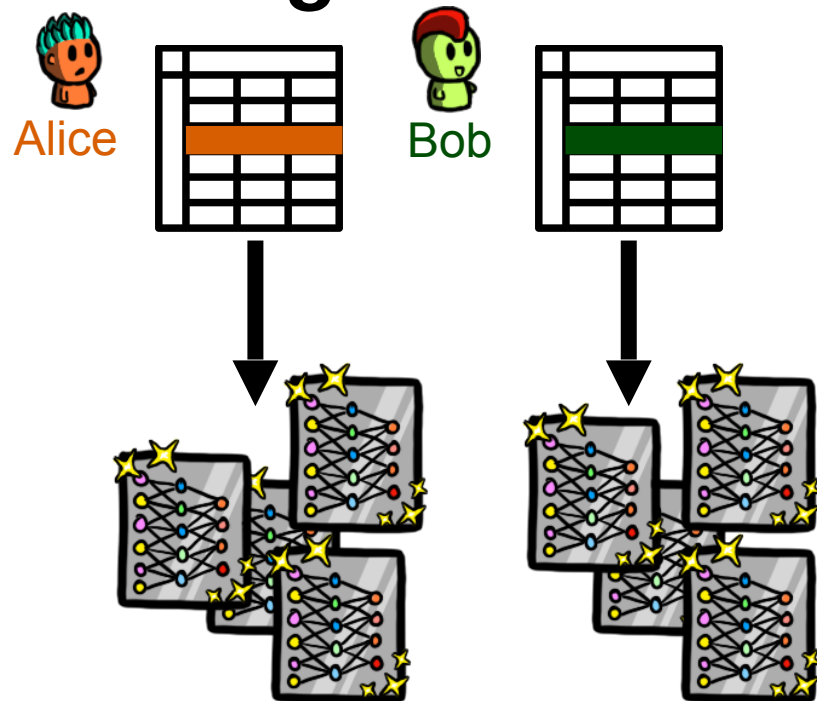
How to make Machine Learning Private

Differential Privacy (Defn.)

Consider any

- Neighbouring datasets S_1 and S_2
- Output set Q

Then Algorithm is (ϵ, δ) -DP if



How to make Machine Learning Private

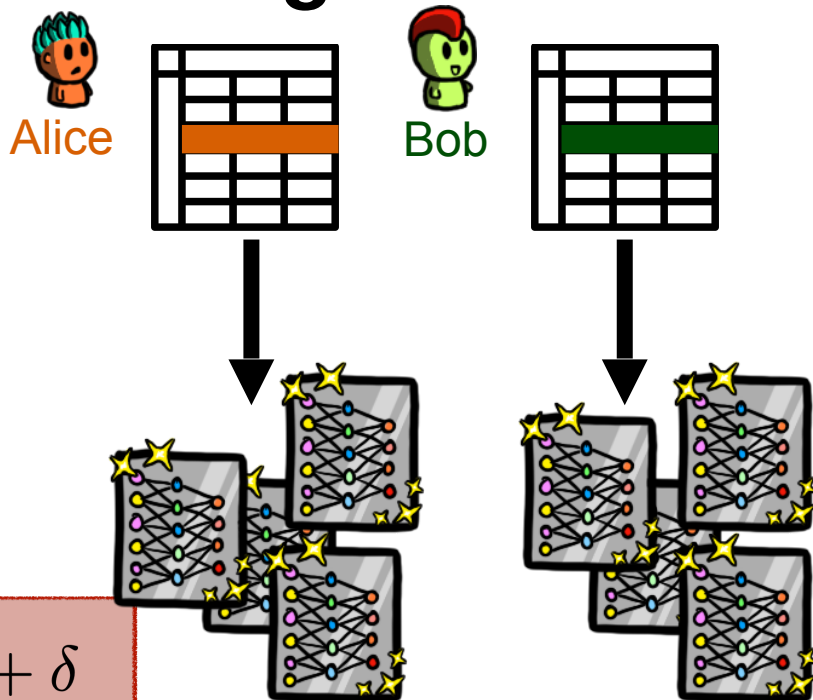
Differential Privacy (Defn.)

Consider any

- Neighbouring datasets S_1 and S_2
- Output set Q

Then Algorithm is (ϵ, δ) -DP if

$$\mathbb{P}(\mathcal{A}(S_1) \in Q) \leq e^\epsilon \mathbb{P}(\mathcal{A}(S_2) \in Q) + \delta$$



Differential Privacy and data quality

Differential Privacy and data quality

One simple formula for implementing DP :

Differential Privacy and data quality

One simple formula for implementing DP :

- **Compute the algorithm's sensitivity:** how much can the output change if the *worst* data point in the *worst* input dataset changes

Differential Privacy and data quality

One simple formula for implementing DP :

- **Compute the algorithm's sensitivity:** how much can the output change if the *worst* data point in the *worst* input dataset changes
- **Add noise** proportional to that change magnitude

Differential Privacy and data quality

One simple formula for implementing DP :

- **Compute the algorithm's sensitivity:** how much can the output change if the *worst* data point in the *worst* input dataset changes
- **Add noise** proportional to that change magnitude

Today, we will look at two ways in which data quality affects the performance of Differentially Private Algorithms

Differential Privacy and data quality

One simple formula for implementing DP :

- **Compute the algorithm's sensitivity:** how much can the output change if the *worst* data point in the *worst* input dataset changes
- **Add noise** proportional to that change magnitude

Today, we will look at two ways in which data quality affects the performance of Differentially Private Algorithms

- *Good data requires less added noise* for the same level of privacy

Differential Privacy and data quality

One simple formula for implementing DP :

- **Compute the algorithm's sensitivity:** how much can the output change if the *worst* data point in the *worst* input dataset changes
- **Add noise** proportional to that change magnitude

Today, we will look at two ways in which data quality affects the performance of Differentially Private Algorithms

- Good data requires less added noise for the same level of privacy
- Some parts of data domain incurs disproportionately higher loss due to the Differential privacy than others

Differential Privacy and Disparate Impact

DP and Disparate Impact

Examples in Practice

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay²,
Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein²,
Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Triedman¹

¹Wikimedia Foundation – htriedman@wikimedia.org

²Tumult Labs – science@tmlt.io

DP and Disparate Impact

Examples in Practice

- Wikimedia foundation released their pageview statistics Differentially Privately.

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay²,
Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein²,
Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Triedman¹

¹Wikimedia Foundation – htriedman@wikimedia.org

²Tumult Labs – science@tmlt.io

DP and Disparate Impact

Examples in Practice

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay²,
Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein²,
Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Triedman¹

¹Wikimedia Foundation – htriedman@wikimedia.org

²Tumult Labs – science@tmlt.io

DP and Disparate Impact

Examples in Practice

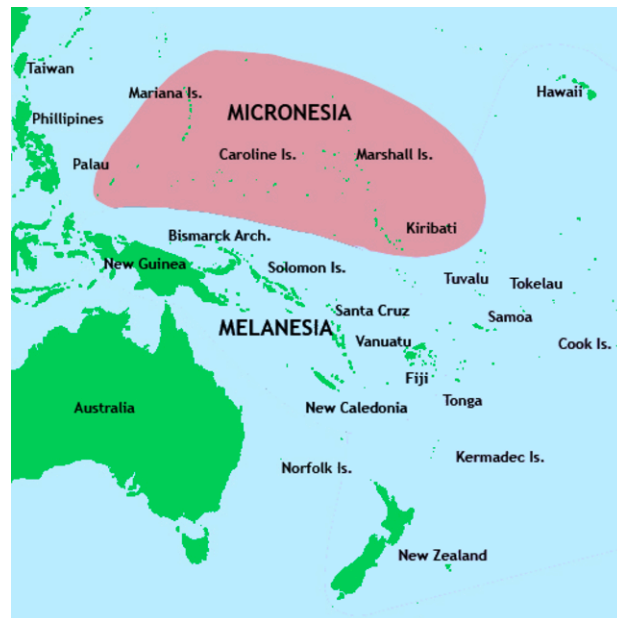
- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay²,
Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein²,
Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Tiedman¹

¹Wikimedia Foundation – htiedman@wikimedia.org

²Tumult Labs – science@tmlt.io



DP and Disparate Impact

Examples in Practice

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

But, the “Micronesia problem”

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Tiedman¹

¹Wikimedia Foundation – htiedman@wikimedia.org

²Tumult Labs – science@tmlt.io



DP and Disparate Impact

Examples in Practice

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

But, the “Micronesia problem”

- Seven Pacific Island nations (low traffic)

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Triedman¹

¹Wikimedia Foundation – htriedman@wikimedia.org

²Tumult Labs – science@tmlt.io



DP and Disparate Impact

Examples in Practice

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

But, the “Micronesia problem”

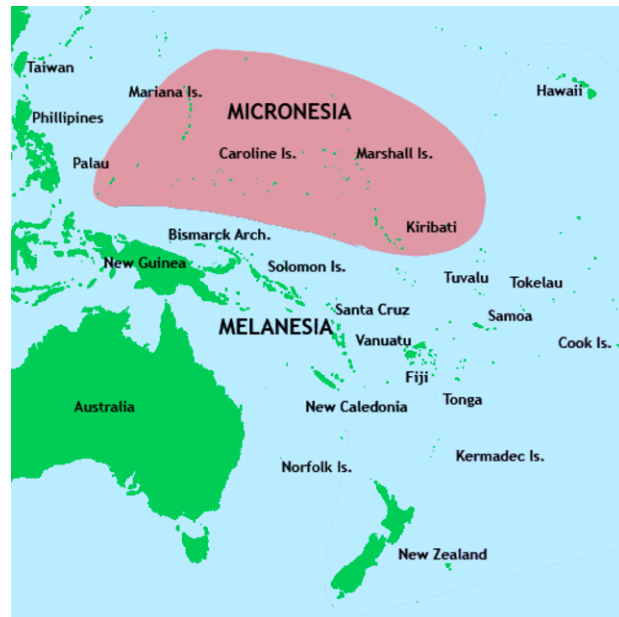
- Seven Pacific Island nations (low traffic)
- Naive first implementation

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Triedman¹

¹Wikimedia Foundation – htriedman@wikimedia.org

²Tumult Labs – science@tmlt.io



DP and Disparate Impact

Examples in Practice

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

But, the “Micronesia problem”

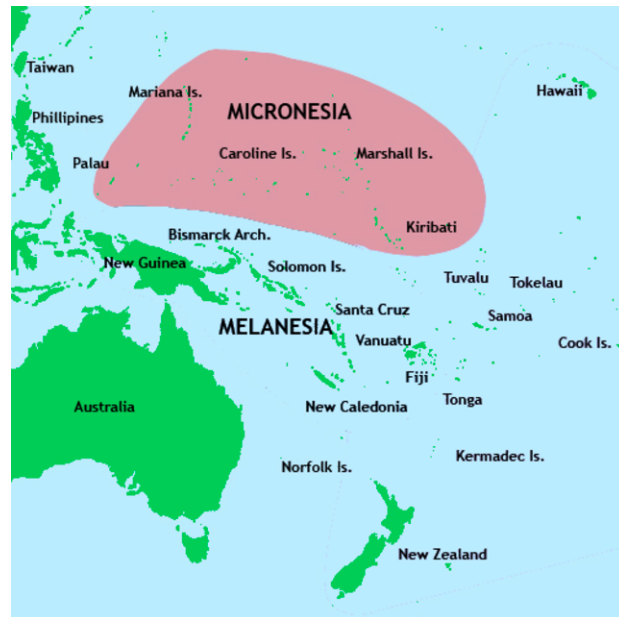
- Seven Pacific Island nations (low traffic)
- Naive first implementation
 - >99% of published data is spurious

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajjhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Tiedman¹

¹Wikimedia Foundation – htiedman@wikimedia.org

²Tumult Labs – science@tmlt.io



DP and Disparate Impact

Examples in Practice

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Tiedman¹

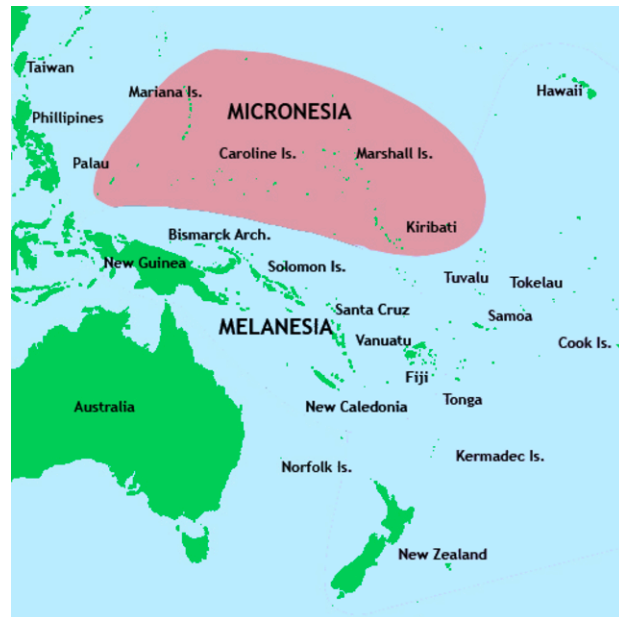
¹Wikimedia Foundation – htiedman@wikimedia.org

²Tumult Labs – science@tmlt.io

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

But, the “Micronesia problem”

- Seven Pacific Island nations (low traffic)
- Naive first implementation
 - >99% of published data is spurious
 - **9 out of 23 subcontinental regions have spurious rate of >25%**



DP and Disparate Impact

Examples in Practice

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Triedman¹

¹Wikimedia Foundation – htriedman@wikimedia.org

²Tumult Labs – science@tmlt.io

- Wikimedia foundation released their pageview statistics Differentially Privately.
- In their global release, both drop rate and spurious rate were less than 1%

But, the “Micronesia problem”

- Seven Pacific Island nations (low traffic)
- Naive first implementation
 - >99% of published data is spurious
 - **9 out of 23 subcontinental regions have spurious rate of >25%**
 - Africa, Oceania, Central Asia, and the Caribbean



DP and Disparate Impact

Examples in Practice

Publishing Wikipedia usage data with strong privacy guarantees

Temilola Adeleye¹, Skye Berghel², Damien Desfontaines², Michael Hay², Isaac Johnson¹, Cléo Lemoisson¹, Ashwin Machanavajhala², Tom Magerlein², Gabriele Modena¹, David Pujol², Daniel Simmons-Marengo², and Hal Tiedman¹

¹Wikimedia Foundation – htiedman@wikimedia.org

²Tumult Labs – science@tmlt.io

- Wikimedia foundation released their pageview statistics Differentially Privately.



The New York Times

The 2020 Census Suggests That People Live Underwater. There's a Reason.

But, the

- Se
- Na
- 9 out of 23 subcontinental regions have spurious rate of >25%
- Africa, Oceania, Central Asia, and the Caribbean

DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

DP and Disparate Impact

Controlled experimental setting

40 binary attributes
for each image

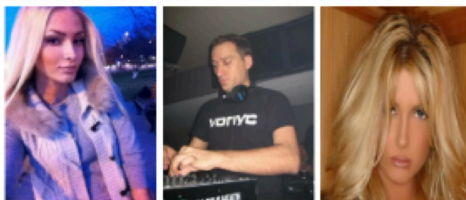
eyeglass



bangs



Pointy nose



How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

DP and Disparate Impact

Controlled experimental setting

40 binary attributes
for each image

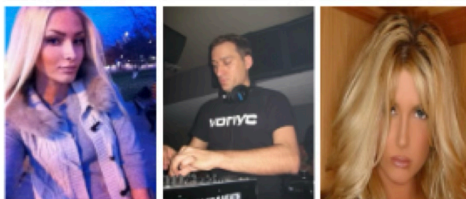
eyeglass



bangs



Pointy nose



40 binary attributes -> 2⁴⁰ subpopulations.

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

DP and Disparate Impact

Controlled experimental setting

40 binary attributes
for each image

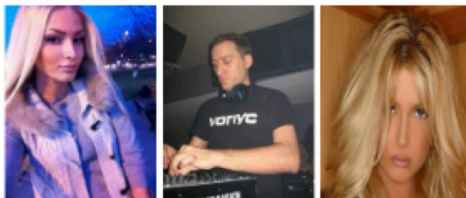
eyeglass



bangs



Pointy nose



40 binary attributes -> 2⁴⁰ subpopulations.

- Subpopulation 1 Eyeglass, no bangs, no pointy nose,...

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

DP and Disparate Impact

Controlled experimental setting

40 binary attributes
for each image

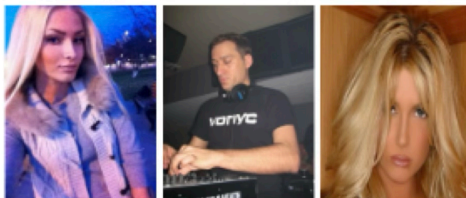
eyeglass



bangs



Pointy nose



40 binary attributes -> 2⁴⁰ subpopulations.

- Subpopulation 1 Eyeglass, no bangs, no pointy nose,...
- Subpopulation 2 Eyeglass, bangs, no pointy nose,...

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

40 binary attributes
for each image

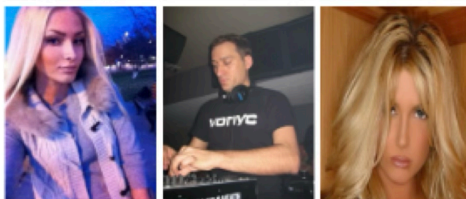
eyeglass



bangs



Pointy nose



40 binary attributes -> 2⁴⁰ subpopulations.

- Subpopulation 1 Eyeglass, no bangs, no pointy nose,...
- Subpopulation 2 Eyeglass, bangs, no pointy nose,...
- ...



DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

40 binary attributes
for each image

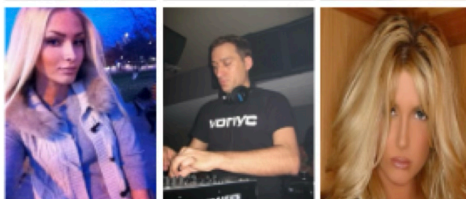
eyeglass



bangs



Pointy nose



40 binary attributes -> 2⁴⁰ subpopulations.

- Subpopulation 1 Eyeglass, no bangs, no pointy nose,...
- Subpopulation 2 Eyeglass, bangs, no pointy nose,...
- ...
- ...



DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu²

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

40 binary attributes
for each image

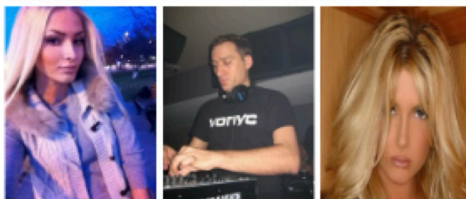
eyeglass



bangs



Pointy nose



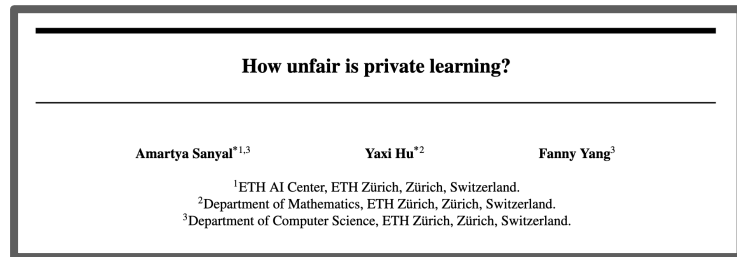
40 binary attributes -> 2⁴⁰ subpopulations.

- Subpopulation 1 Eyeglass, no bangs, no pointy nose,...
- Subpopulation 2 Eyeglass, bangs, no pointy nose,...
- ...
- ...
- Subpopulation 2⁴⁰



DP and Disparate Impact

Controlled experimental setting



40 binary attributes
for each image

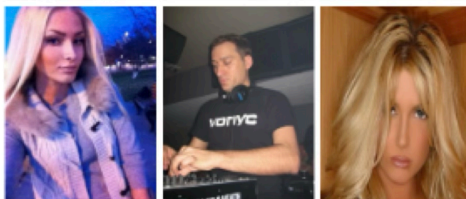
eyeglass



bangs

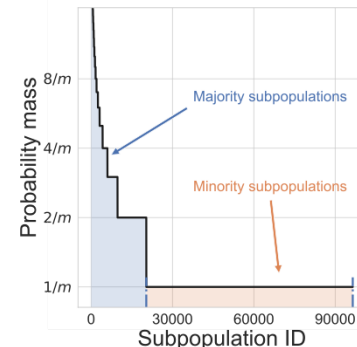


Pointy nose



40 binary attributes -> 2⁴⁰ subpopulations.

- Subpopulation 1 Eyeglass, no bangs, no pointy nose,...
- Subpopulation 2 Eyeglass, bangs, no pointy nose,...
- ...
- ...
- Subpopulation 2⁴⁰



DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

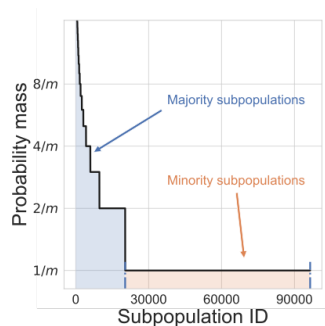
Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.



DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

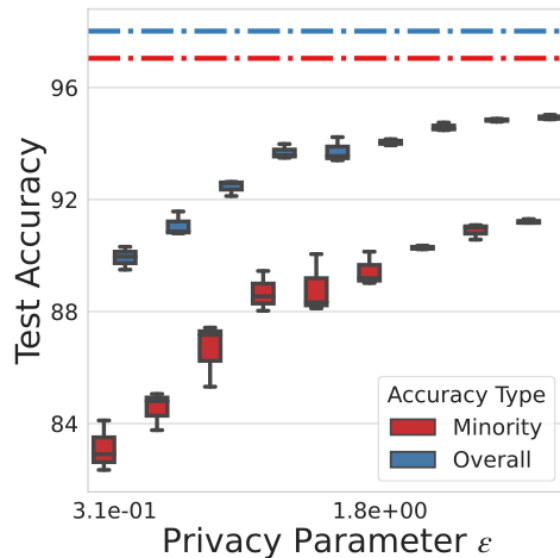
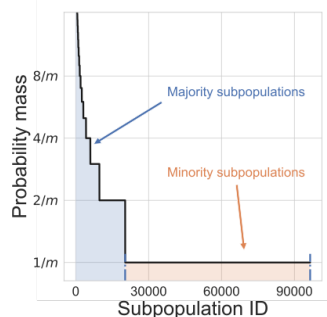
Yaxi Hu²

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.



DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

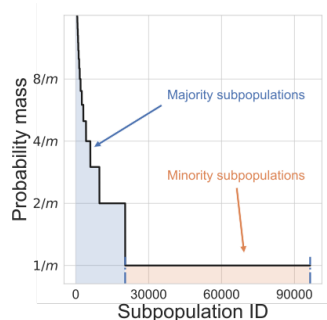
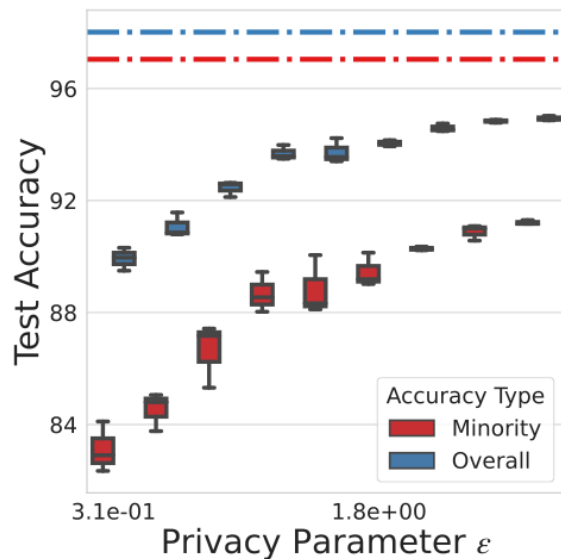
Yaxi Hu²

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.



DP and Disparate Impact

Controlled experimental setting

How unfair is private learning?

Amartya Sanyal^{1,3}

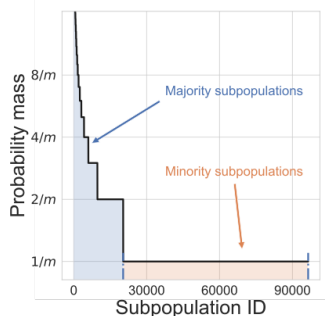
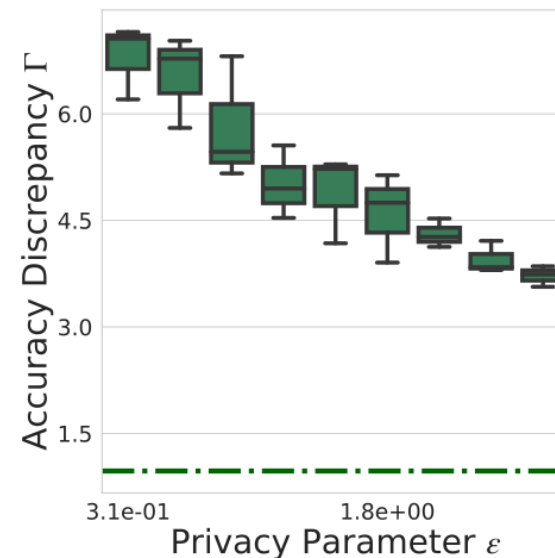
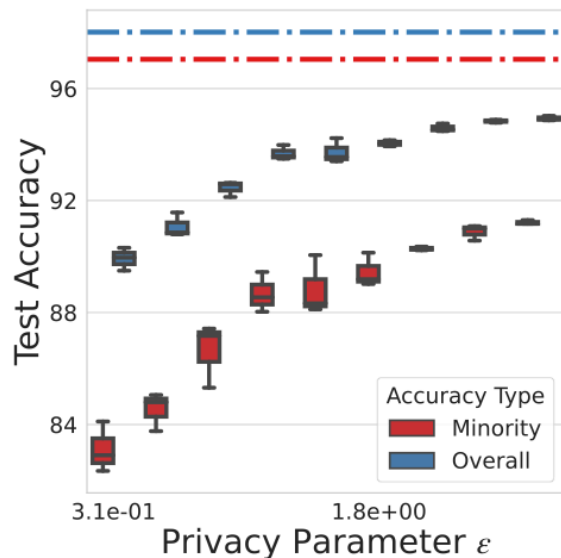
Yaxi Hu²

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.

DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu²

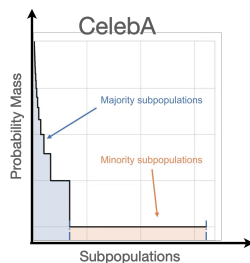
Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu²

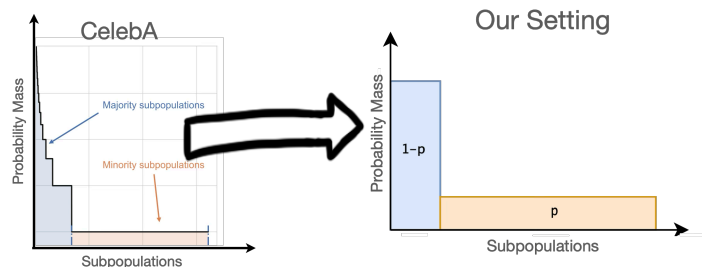
Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu²

Fanny Yang³

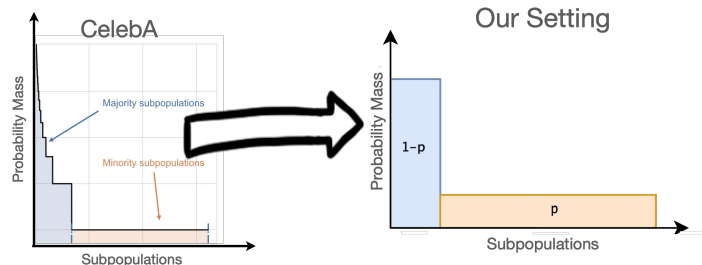
¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.

$$\text{If } \frac{N}{m} \rightarrow c \text{ as } N, m \rightarrow \infty$$



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu²

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

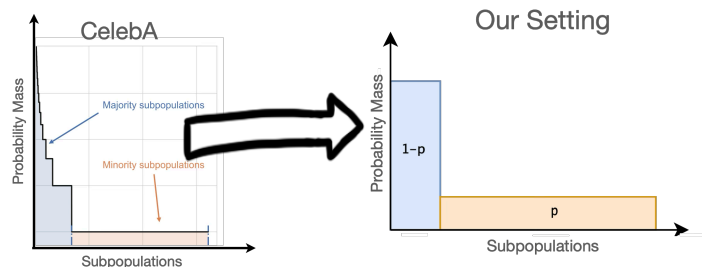
³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.

$$\text{If } \frac{N}{m} \rightarrow c \text{ as } N, m \rightarrow \infty$$

We prove a lower bound on the *accuracy discrepancy* of \mathcal{A} which

- **Increases with the privacy** of \mathcal{A} i.e. smaller (ϵ, δ)



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

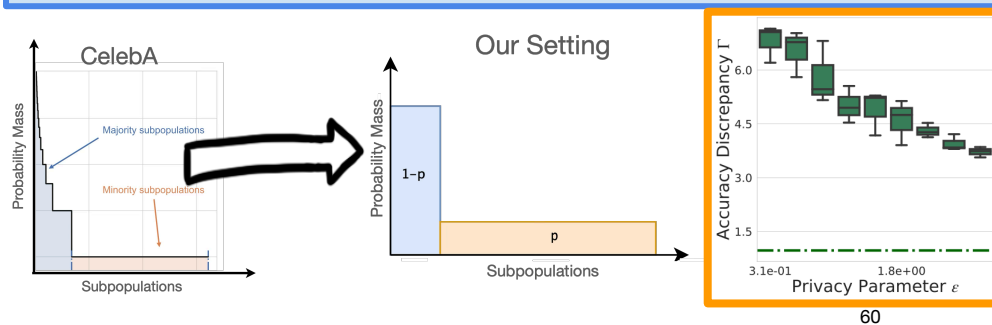
³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.

$$\text{If } \frac{N}{m} \rightarrow c \text{ as } N, m \rightarrow \infty$$

We prove a lower bound on the *accuracy discrepancy* of \mathcal{A} which

- Increases with the privacy of \mathcal{A} i.e. smaller (ϵ, δ)



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

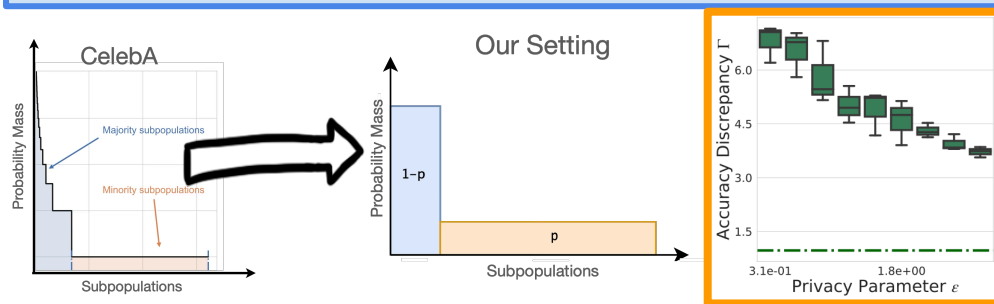
³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.

$$\text{If } \frac{N}{m} \rightarrow c \text{ as } N, m \rightarrow \infty$$

We prove a lower bound on the *accuracy discrepancy* of \mathcal{A} which

- Increases with the privacy of \mathcal{A} i.e. smaller (ϵ, δ)
- Increases with **long-tailed** nature of data i.e. relative number of minority subpopulations c



DP and Disparate Impact

Trade-off in long-tailed data

How unfair is private learning?

Amartya Sanyal^{1,3}

Yaxi Hu^{*2}

Fanny Yang³

¹ETH AI Center, ETH Zürich, Zürich, Switzerland.

²Department of Mathematics, ETH Zürich, Zürich, Switzerland.

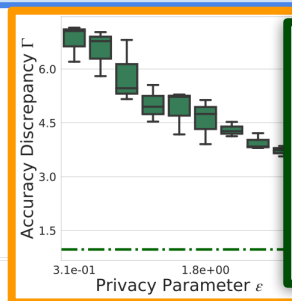
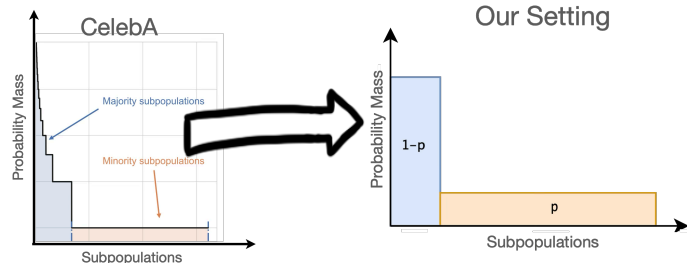
³Department of Computer Science, ETH Zürich, Zürich, Switzerland.

Theorem Consider a *long-tailed distribution* with N sub-populations and sample an m -sized dataset. Consider any (ϵ, δ) -DP algorithm \mathcal{A} that achieves *low error* on this dataset.

$$\text{If } \frac{N}{m} \rightarrow c \text{ as } N, m \rightarrow \infty$$

We prove a lower bound on the *accuracy discrepancy* of \mathcal{A} which

- Increases with the privacy of \mathcal{A} i.e. smaller (ϵ, δ)
- Increases with **long-tailed** nature of data i.e. relative number of minority subpopulations c



DP and Disparate Impact

Fundamental Impossibility

DP and Disparate Impact

Fundamental Impossibility

Trade-Offs between Fairness and Privacy in Machine Learning

Sushant Agarwal
University of Waterloo, Canada
sushant.agarwal@uwaterloo.ca

On the Compatibility of Privacy and Fairness

Rachel Cummings* Varun Gupta* Dhamma Kimpara* Jamie Morgenstern*

DP and Disparate Impact

Fundamental Impossibility

Trade-Offs between Fairness and Privacy in Machine Learning

Sushant Agarwal
University of Waterloo, Canada
sushant.agarwal@uwaterloo.ca

On the Compatibility of Privacy and Fairness

Rachel Cummings* Varun Gupta* Dhamma Kimpara* Jamie Morgenstern*

Theorem For any hypothesis class \mathcal{H} , no algorithm can simultaneously be $(\epsilon, 0)$ -DP for $\epsilon < \infty$ and always output a $h \in \mathcal{H}$ that satisfies *equal opportunity* and has *error less than for any constant classifier*.

DP and Disparate Impact

Fundamental Impossibility

Trade-Offs between Fairness and Privacy in Machine Learning

Sushant Agarwal
University of Waterloo, Canada
sushant.agarwal@uwaterloo.ca

On the Compatibility of Privacy and Fairness

Rachel Cummings* Varun Gupta* Dhamma Kimpara* Jamie Morgenstern*

Theorem For any hypothesis class \mathcal{H} , no algorithm can simultaneously be $(\epsilon, 0)$ -DP for $\epsilon < \infty$ and always output a $h \in \mathcal{H}$ that satisfies *equal opportunity* and has *error less than for any constant classifier*.

- Proof idea:
 - Obs 1. - If a classifier h has non-zero probability to be output under algorithm \mathcal{A} on dataset S_1 , it also has non-zero probability to be output on dataset S_2 , for all S_2 .

DP and Disparate Impact

Fundamental Impossibility

Trade-Offs between Fairness and Privacy in Machine Learning

Sushant Agarwal
University of Waterloo, Canada
sushant.agarwal@uwaterloo.ca

On the Compatibility of Privacy and Fairness

Rachel Cummings* Varun Gupta* Dhamma Kimpara* Jamie Morgenstern*

Theorem For any hypothesis class \mathcal{H} , no algorithm can simultaneously be $(\epsilon, 0)$ -DP for $\epsilon < \infty$ and always output a $h \in \mathcal{H}$ that satisfies *equal opportunity* and has *error less than for any constant classifier*.

- Proof idea:
 - Obs 1. - If a classifier h has non-zero probability to be output under algorithm \mathcal{A} on dataset S_1 , it also has non-zero probability to be output on dataset S_2 , for all S_2 .
 - Obs 2. - Construct two datasets S_1, S_2 such that no classifier, except a constant classifier can be simultaneously fair on both.

DP and Disparate Impact

Other causes

DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

In particular two factors are usually isolated,

DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

In particular two factors are usually isolated,

- **Gradient Clipping**

DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

In particular two factors are usually isolated,

- **Gradient Clipping**
- **Noise** addition

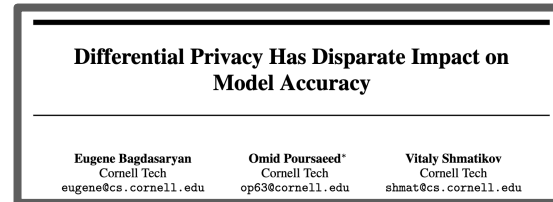
DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

In particular two factors are usually isolated,

- **Gradient Clipping**
- **Noise** addition



DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

In particular two factors are usually isolated,

- **Gradient Clipping**
- **Noise** addition

Differentially Private Empirical Risk Minimization under the Fairness Lens		
Cuong Tran Syracuse University ctran@syr.edu	My H. Dinh Syracuse University mydinh@syr.edu	Ferdinando Fioretto Syracuse University ffiorett@syr.edu

Differential Privacy Has Disparate Impact on Model Accuracy		
Eugene Bagdasaryan Cornell Tech eugene@cs.cornell.edu	Omid Poursaeed* Cornell Tech op63@cornell.edu	Vitaly Shmatikov Cornell Tech shmat@cs.cornell.edu

DP and Disparate Impact

Other causes

Apart from the properties of the data, other reasons are also known to exacerbate unfairness for Differentially Private models

In particular two factors are usually isolated,

- **Gradient Clipping**
- **Noise** addition

Differentially Private Empirical Risk Minimization under the Fairness Lens

Cuong Tran
Syracuse University
ctran@syr.edu

My H. Dinh
Syracuse University
mydinh@syr.edu

Ferdinando Fioretto
Syracuse University
ffiorett@syr.edu

Differential Privacy Has Disparate Impact on Model Accuracy

Eugene Bagdasaryan
Cornell Tech
eugene@cs.cornell.edu

Omid Poursaeed*
Cornell Tech
op63@cornell.edu

Vitaly Shmatikov
Cornell Tech
shmat@cs.cornell.edu

Removing Disparate Impact on Model Accuracy in Differentially Private Stochastic Gradient Descent

Depeng Xu
University of Arkansas
Fayetteville, AR, USA
depengxu@uark.edu

Wei Du
University of Arkansas
Fayetteville, AR, USA
wd005@uark.edu

Xintao Wu
University of Arkansas
Fayetteville, AR, USA
xintaowu@uark.edu

Good data requires less noise

DP with “good” data

Favourable data properties

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.
- Generally, these worst case bound yields a large “*cost of DP*”.

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.
- Generally, these worst case bound yields a large “*cost of DP*”.

Differentially Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds

Raef Bassily*

Adam Smith*[†]

Abhradeep Thakurta[‡]

Theorem For every convex function in \mathbb{R}^d , DP-SGD satisfies DP and cost of DP is less than $\frac{LR\sqrt{d}}{\epsilon n}$.

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.
- Generally, these worst case bound yields a large “*cost of DP*”.

Differentially Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds

Raef Bassily*

Adam Smith*[†]

Abhradeep Thakurta[‡]

Theorem For every convex function in \mathbb{R}^d , DP-SGD satisfies DP and cost of DP is less than $\frac{LR\sqrt{d}}{\epsilon n}$.

Theorem Exists a convex loss function and data in \mathbb{R}^d , s.t. all DP algorithms incur extra cost $\frac{LR\sqrt{d}}{\epsilon n}$.

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.
- Generally, these worst case bound yields a large “cost of DP”.

Can we do better for “nice” datasets ?

Theorem For every convex function in \mathbb{R}^d , DP-SGD satisfies DP and cost of DP is less than $\frac{LR\sqrt{d}}{\epsilon n}$.

Theorem Exists a convex loss function and data in \mathbb{R}^d , s.t. all DP algorithms incur extra cost $\frac{LR\sqrt{d}}{\epsilon n}$.

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.
- Generally, these worst case bound yields a large “cost of DP”.

Can we do better for “nice” datasets ?

Differentially Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds

Raef Bassily*

Adam Smith*†

Abhradeep Thakurta‡



Theorem For every convex function in \mathbb{R}^d , DP-SGD satisfies DP and cost of DP is less than $\frac{LR\sqrt{d}}{\epsilon n}$.

Theorem Exists a convex loss function and data in \mathbb{R}^d , s.t. all DP algorithms incur extra cost $\frac{LR\sqrt{d}}{\epsilon n}$.

DP with “good” data

Favourable data properties

- Privacy guarantees are unconditional - hold for all datasets.
- Generally, these worst case bound yields a large “cost of DP”.

Can we do better for
“nice” datasets ?

Theorem For every convex function in \mathbb{R}^d , DP-SGD satisfies DP and cost of DP is less than $\frac{LR\sqrt{d}}{\epsilon n}$.

Theorem Exists a convex loss function and data in \mathbb{R}^d , s.t. all DP algorithms incur extra cost $\frac{LR\sqrt{d}}{\epsilon n}$.

Differentially Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds

Raef Bassily*

Adam Smith*†

Abhradeep Thakurta‡



Naive DP Mean estimation

Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;

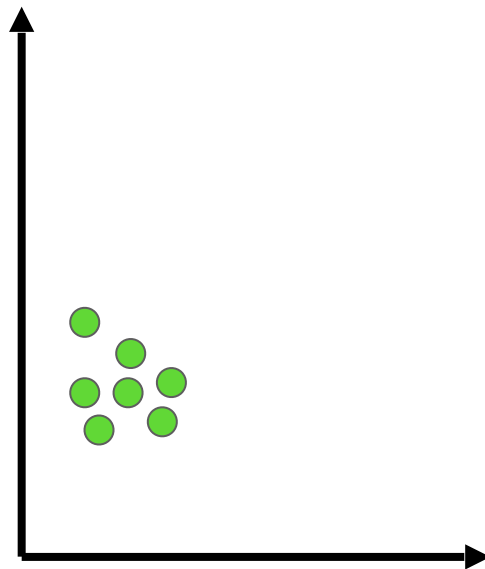
Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;



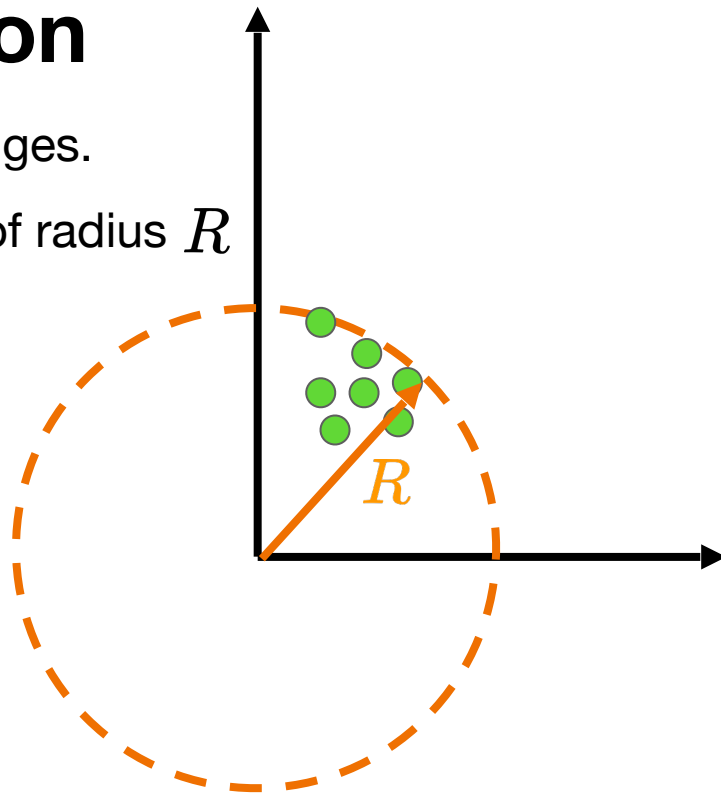
Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;



Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;

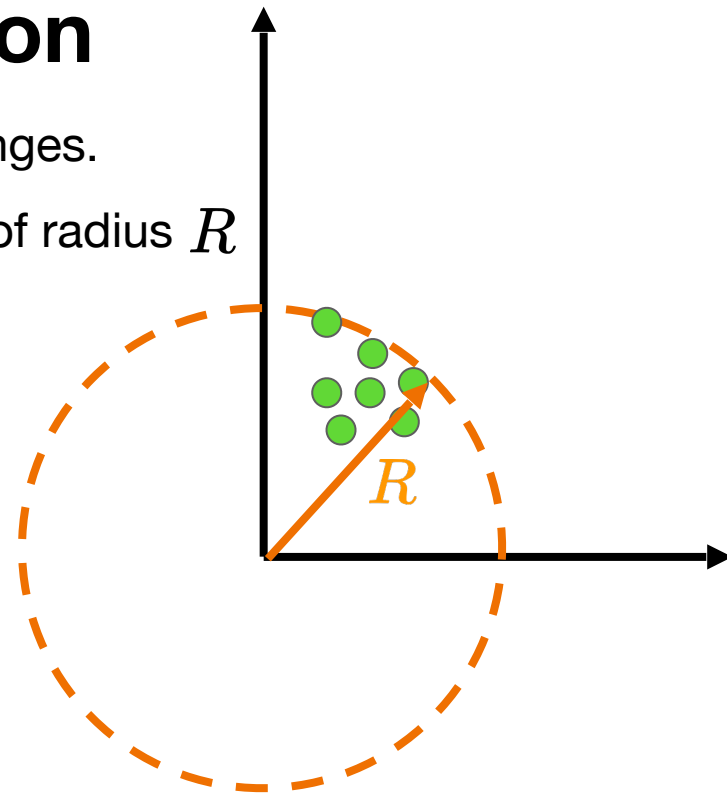


Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;

- error scales as

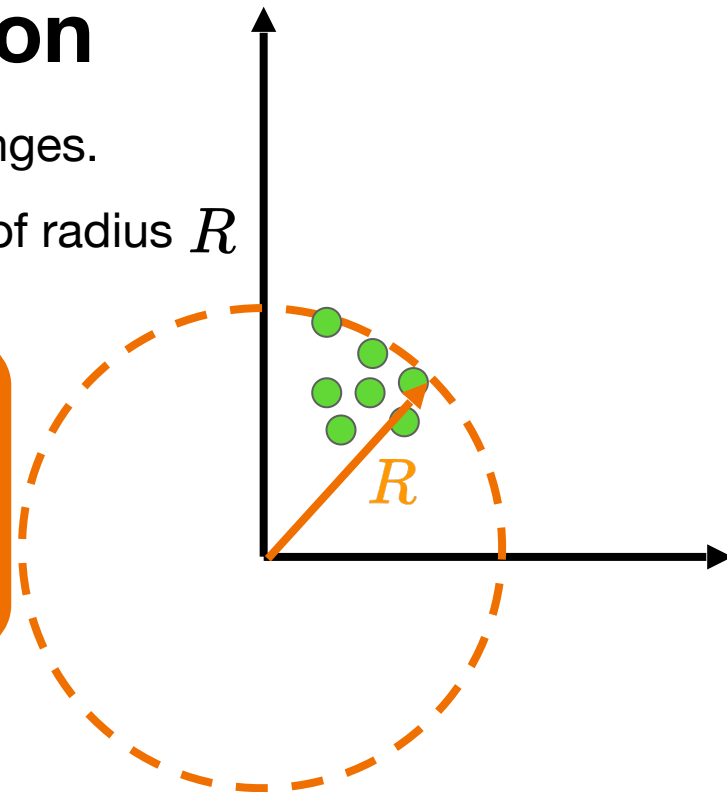
$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n} + \frac{R^2 d^2}{n^2 \epsilon}}$$



Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;
- error scales as

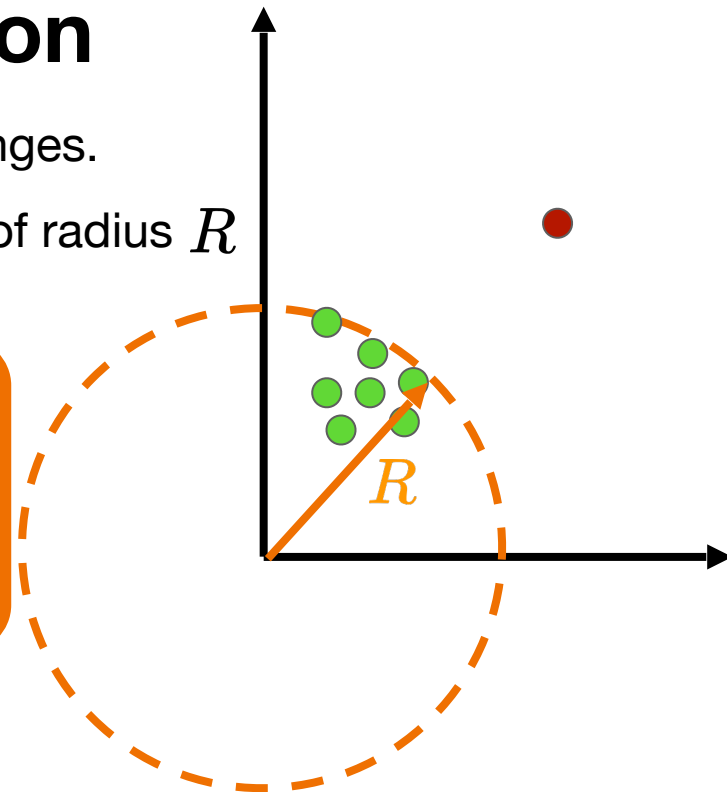
$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n} + \frac{R^2 d^2}{n^2 \epsilon}}$$



Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;
- error scales as

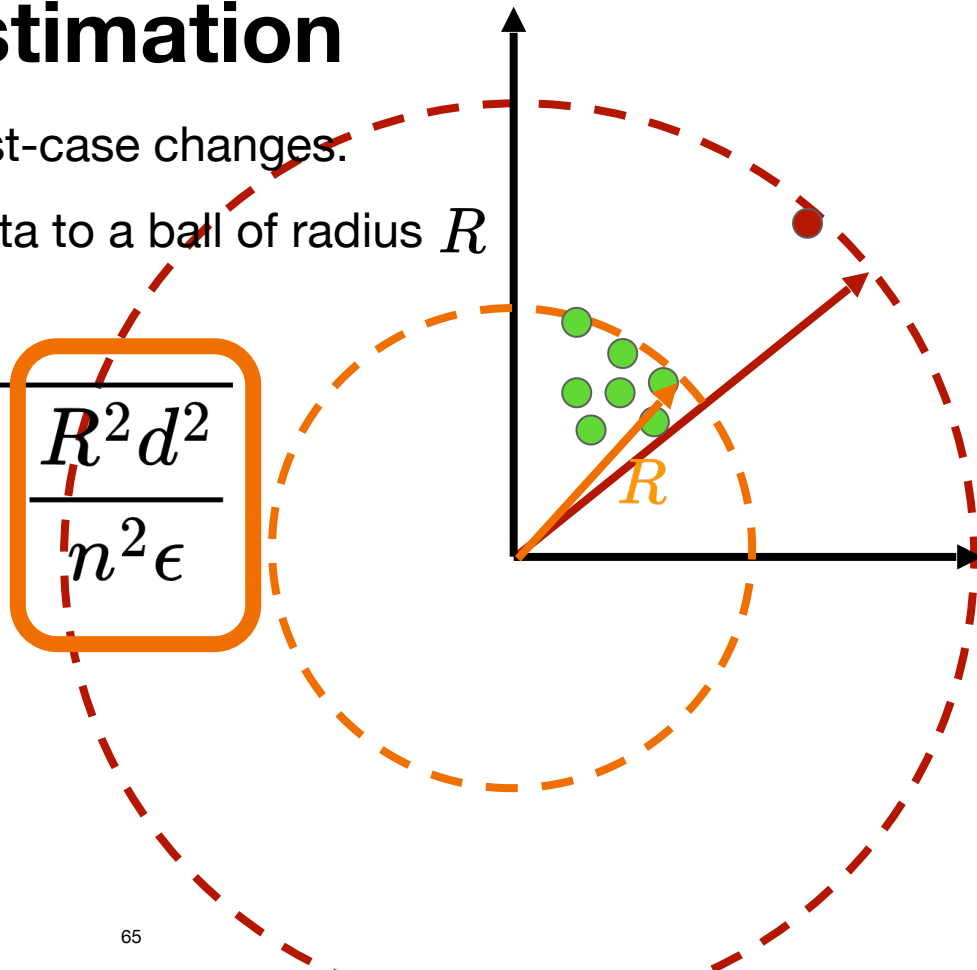
$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n} + \frac{R^2 d^2}{n^2 \epsilon}}$$



Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;
- error scales as

$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n} + \frac{R^2 d^2}{n^2 \epsilon}}$$

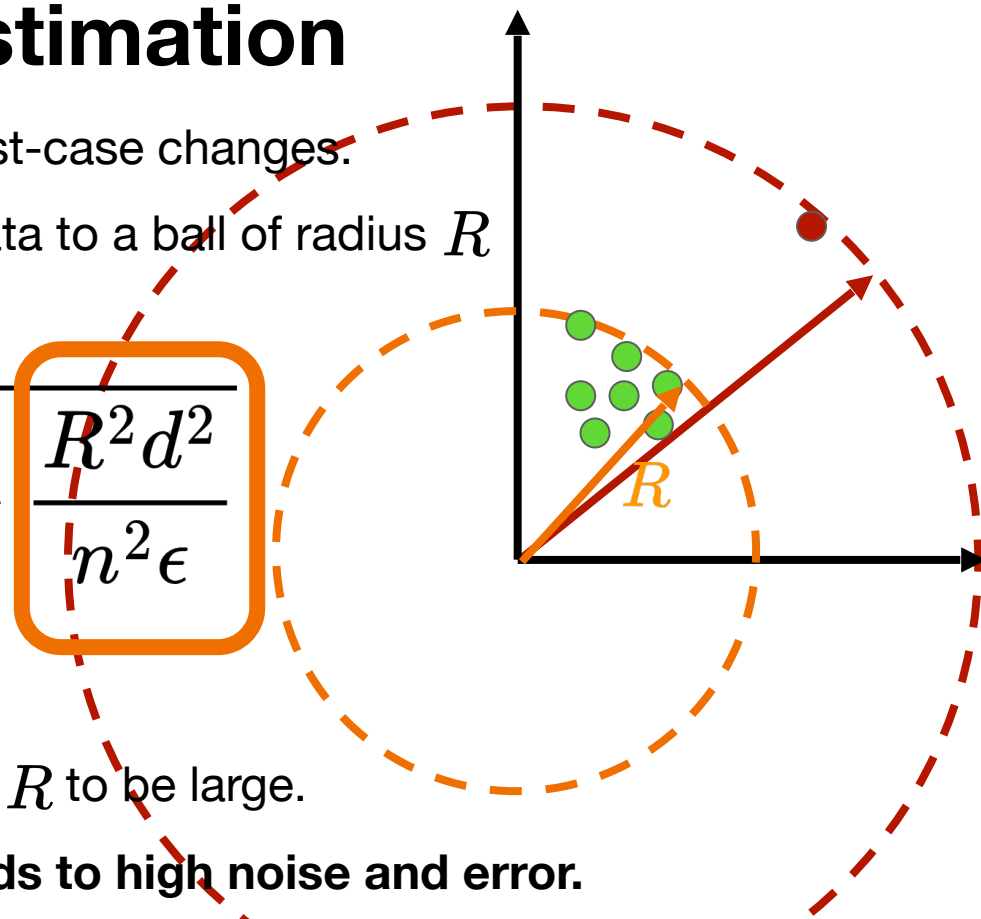


Naive DP Mean estimation

- DP must protect against worst-case changes.
- A naïve DP estimator clips data to a ball of radius R
- Noise is calibrated to R ;

- error scales as

$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n} + \frac{R^2 d^2}{n^2 \epsilon}}$$



- Heavy tails or outliers force R to be large.
- **Worst-case sensitivity leads to high noise and error.**

Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

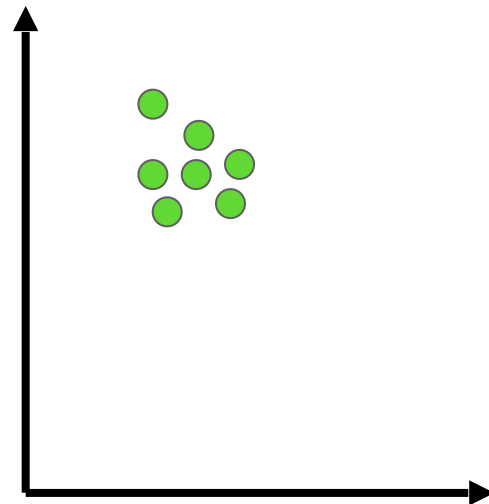
In many settings, most points lie in a ball of radius

Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

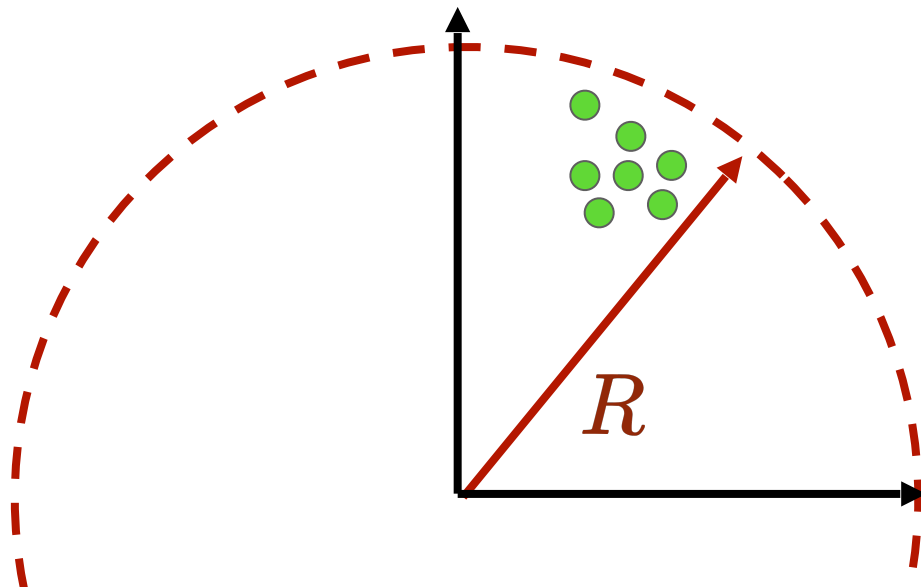


Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius



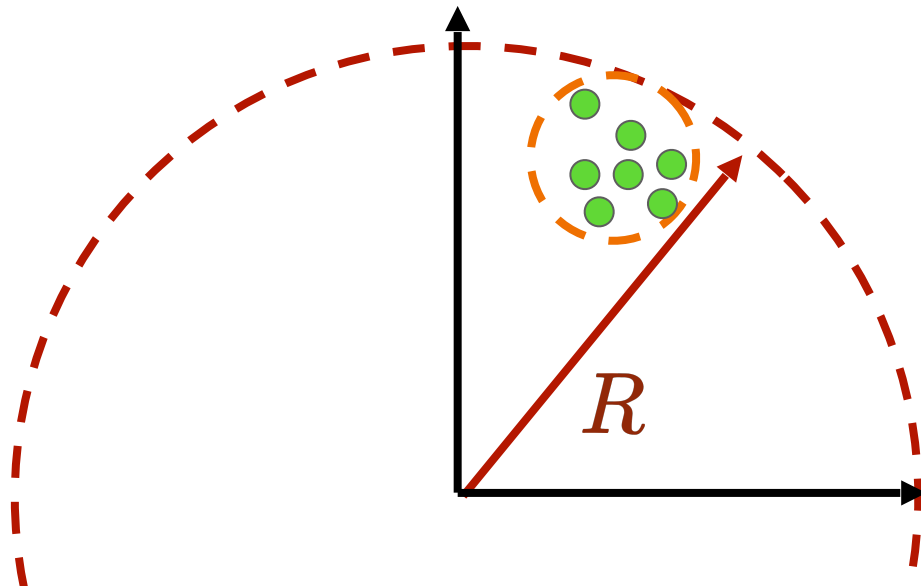
Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

- A dataset is **friendly** if every two points have a common neighbor within r



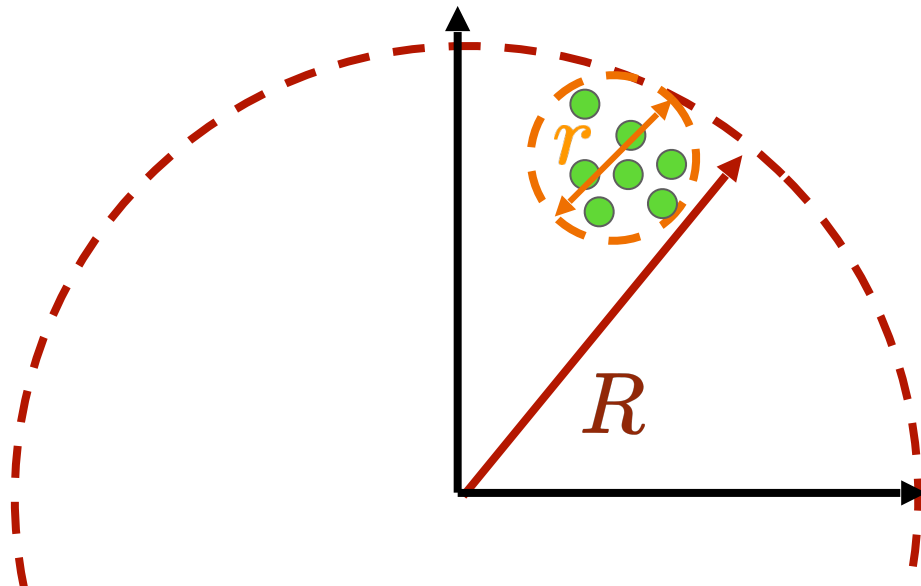
Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

- A dataset is **friendly** if every two points have a common neighbor within r



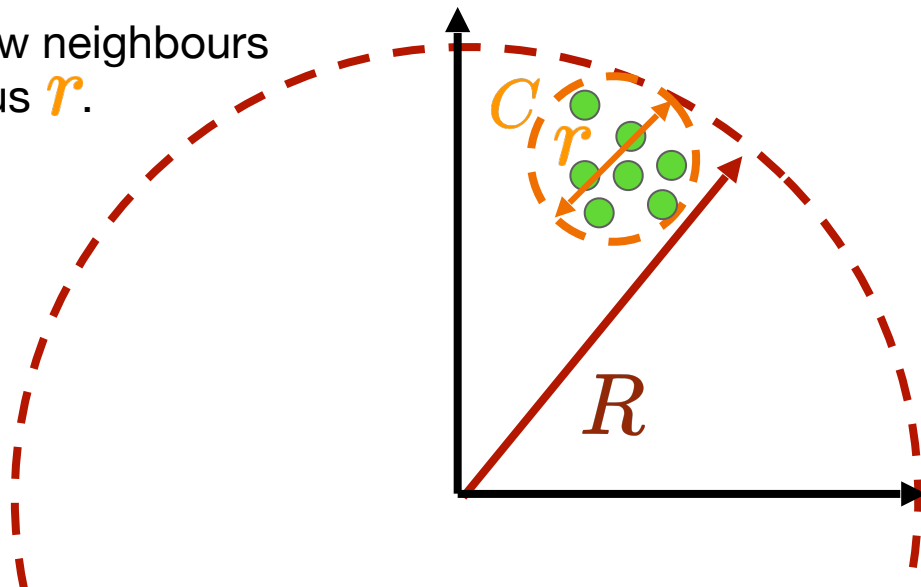
Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

- A dataset is **friendly** if every two points have a common neighbor within r
- **Friendly-Core** removes points with few neighbours and outputs w.h.p. a core C with radius r .



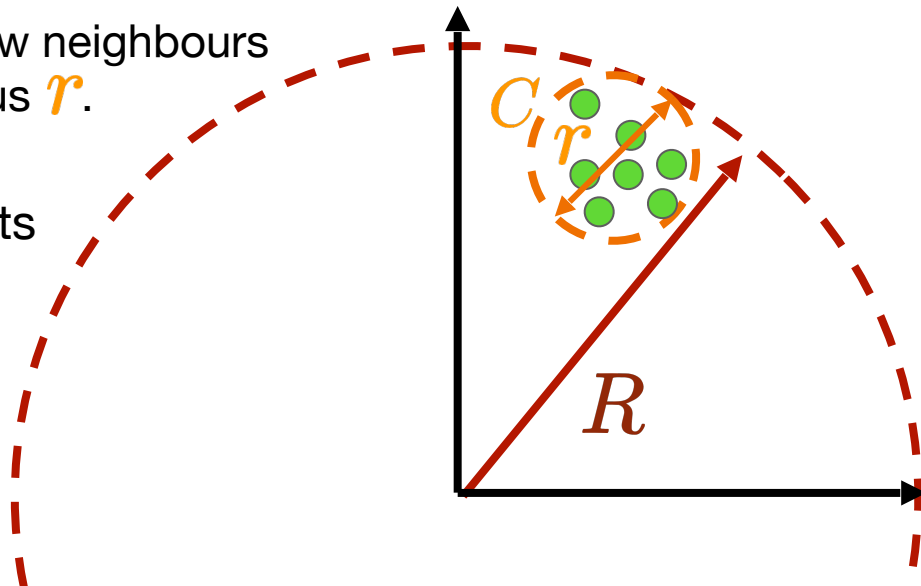
Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

- A dataset is **friendly** if every two points have a common neighbor within r
- **Friendly-Core** removes points with few neighbours and outputs w.h.p. a core \mathcal{C} with radius r .
- Estimating mean in the core \mathcal{C} results



Friendly-Core estimation

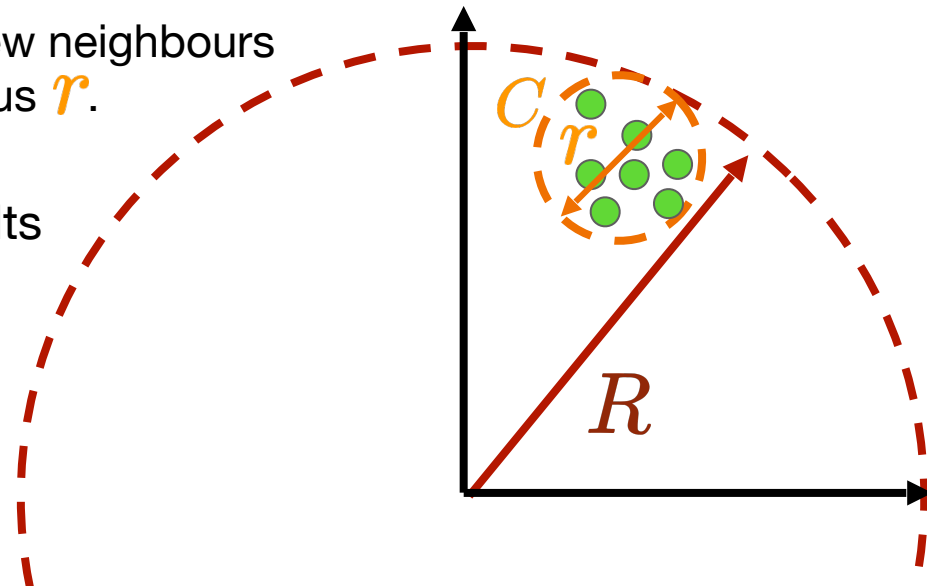
FriendlyCore: Practical Differentially Private Aggregation

Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

- A dataset is **friendly** if every two points have a common neighbor within r
- **Friendly-Core** removes points with few neighbours and outputs w.h.p. a core \mathcal{C} with radius r .
- Estimating mean in the core \mathcal{C} results

$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n}} + \frac{r\sqrt{d}}{n\epsilon}$$



Friendly-Core estimation

FriendlyCore: Practical Differentially Private Aggregation

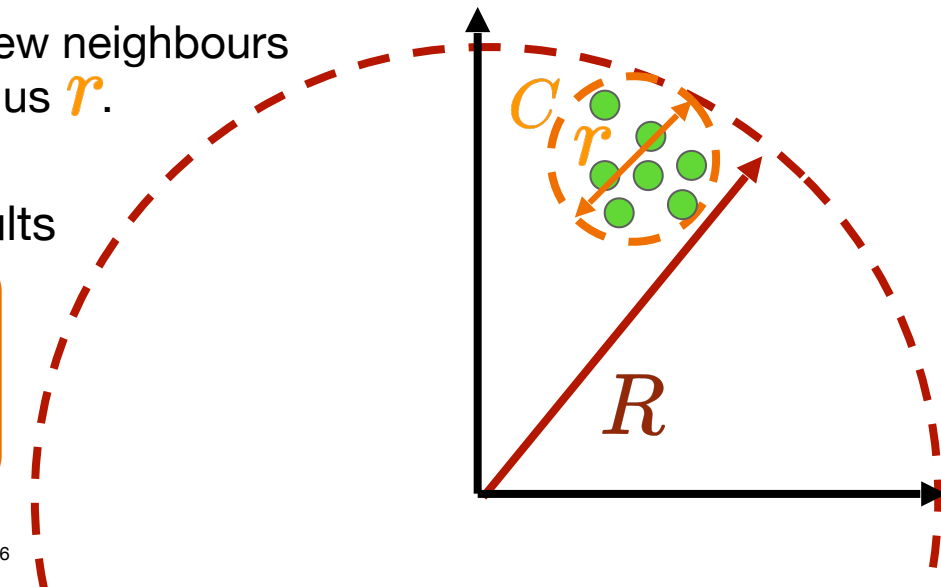
Eliad Tsfadia* Edith Cohen* Haim Kaplan* Yishay Mansour*
Uri Stemmer*

In many settings, most points lie in a ball of radius

- A dataset is **friendly** if every two points have a common neighbor within r
- **Friendly-Core** removes points with few neighbours and outputs w.h.p. a core \mathcal{C} with radius r .
- Estimating mean in the core \mathcal{C} results

$$\|\hat{\mu} - \mu\| \leq \sqrt{\frac{d}{n}} + \frac{r\sqrt{d}}{n\epsilon}$$

when the data is **friendly**.



DP geometric median

Private Geometric Median

Mahdi Haghifam*

Thomas Steinke†

Jonathan Ullman‡

DP geometric median

- Related task is estimating the geometric median: solving the following

$$\sum \| \theta - x_i \|_2$$

DP geometric median

Private Geometric Median

Mahdi Haghifam*

Thomas Steinke†

Jonathan Ullman‡

- Related task is estimating the geometric median: solving the following

$$\sum \| \theta - x_i \|_2$$

- Solving this problem with DP-SGD yields DP cost $\frac{R\sqrt{d}}{\epsilon}$

DP geometric median

- Related task is estimating the geometric median: solving the following

$$\sum \| \theta - x_i \|_2$$

- Solving this problem with DP-SGD yields DP cost $\frac{R\sqrt{d}}{\epsilon}$
- **HSU24** proposes an algorithm which yields DP cost ϵ

DP geometric median

- Related task is estimating the geometric median: solving the following

$$\sum \| \theta - x_i \|_2$$

- Solving this problem with DP-SGD yields DP cost $\frac{R\sqrt{d}}{\epsilon}$
- **HSU24** proposes an algorithm which yields DP cost $(\text{effective diameter}) \frac{\sqrt{d}}{\epsilon}$

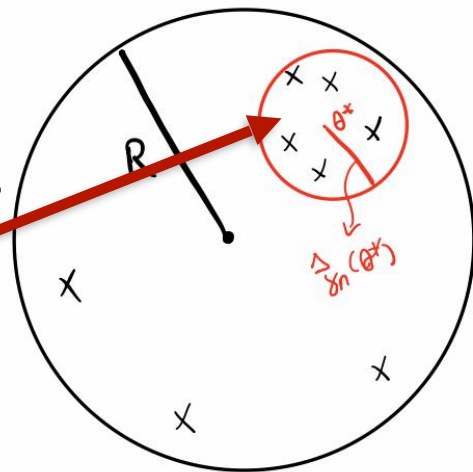
DP geometric median

- Related task is estimating the geometric median: solving the following

$$\sum \|\theta - x_i\|_2$$

- Solving this problem with DP-SGD yields DP cost $\frac{R\sqrt{d}}{\epsilon}$

- **HSU24** proposes an algorithm which yields DP cost $\frac{\epsilon}{\sqrt{d}}$
(effective diameter)



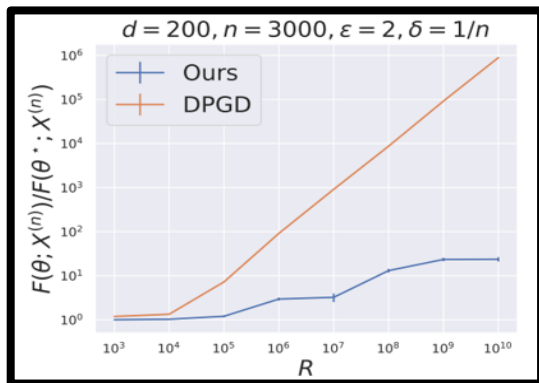
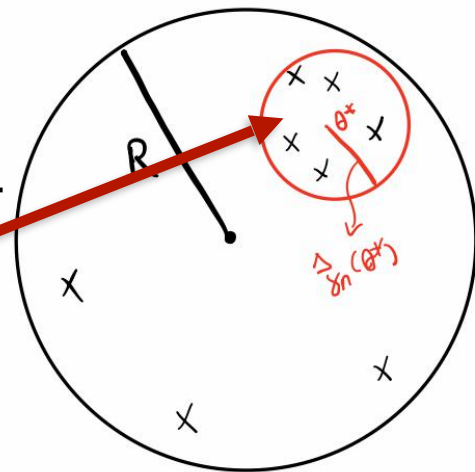
DP geometric median

- Related task is estimating the geometric median: solving the following

$$\sum \| \theta - x_i \|_2$$

- Solving this problem with DP-SGD yields DP cost $\frac{R\sqrt{d}}{\epsilon}$
- **HSU24** proposes an algorithm which yields DP cost $\frac{\epsilon}{\sqrt{d}}$

(effective diameter) $\frac{\sqrt{d}}{\epsilon}$



DP Principal Component Analysis

DP Principal Component Analysis

DP-PCA: Statistically Optimal and
Differentially Private PCA

Xiyang Liu *

Weihaio Kong †

Prateek Jain ‡

Sewoong Oh §

DP Principal Component Analysis

Lower bound for any data

DP-PCA: Statistically Optimal and
Differentially Private PCA

Xiyang Liu * Weihao Kong † Prateek Jain ‡ Sewoong Oh §

Theorem 5.4 (Lower bound without Assumptions A.1–A.3) *that map n i.i.d. samples to an estimate $\hat{v} \in \mathbb{R}^d$. A set of distributions satisfying Assumptions A.1–A.3 with $M = \tilde{O}(d + \sqrt{n\varepsilon/d})$, $V = O(d)$ and $\gamma = O(1)$ is denoted by $\tilde{\mathcal{P}}$. For $d \geq 2$, there exists a universal constant $C > 0$ such that*

$$\inf_{\hat{v} \in \mathcal{M}_\varepsilon} \sup_{P \in \tilde{\mathcal{P}}} \mathbb{E}_{S \sim P^n} [\sin(\hat{v}(S), v_1)] \geq C\kappa \min \left(\sqrt{\frac{d \wedge \log((1 - e^{-\varepsilon})/\delta)}{\varepsilon n}}, 1 \right). \quad (13)$$

DP Principal Component Analysis

Lower bound for any data

DP-PCA: Statistically Optimal and
Differentially Private PCA

Xiyang Liu * Weihaio Kong † Prateek Jain ‡ Sewoong Oh §

Theorem 5.4 (Lower bound without Assumptions A.1–A.3) *Let \mathcal{M}_ϵ be a set of DP-PCA estimators that map n i.i.d. samples to an estimate $\hat{v} \in \mathbb{R}^d$. A set of distributions satisfying Assumptions A.1–A.3 with $M = \tilde{O}(d + \sqrt{n\epsilon/d})$, $V = O(d)$ and $\gamma = O(1)$ is denoted by $\tilde{\mathcal{P}}$. For $d \geq 2$, there exists a universal constant $C > 0$ such that*

$$\inf_{\hat{v} \in \mathcal{M}_\epsilon} \sup_{P \in \tilde{\mathcal{P}}} \mathbb{E}_{S \sim P^n} [\sin(\hat{v}(S), v_1)] \geq C \kappa \min \left(\sqrt{\frac{d \wedge \log((1 - e^{-\epsilon})/\delta)}{\epsilon n}}, 1 \right). \quad (13)$$

Theorem 5.5 (Lower bound, Gaussian distribution) *Let \mathcal{M}_ϵ be a set of DP-PCA estimators that map n i.i.d. samples to an estimate $\hat{v} \in \mathbb{R}^d$. A set of Gaussian distributions with (λ_1, λ_2) as the first and second eigenvalues of the covariance matrix is denoted by $\mathcal{P}_{(\lambda_1, \lambda_2)}$. There exists a universal constant $C > 0$ such that*

$$\inf_{\hat{v} \in \mathcal{M}_\epsilon} \sup_{P \in \mathcal{P}_{(\lambda_1, \lambda_2)}} \mathbb{E}_{S \sim P^n} [\sin(\hat{v}(S), v_1)] \geq C \min \left(\kappa \left(\sqrt{\frac{d}{n}} + \frac{d}{\epsilon n} \right) \sqrt{\frac{\lambda_2}{\lambda_1}}, 1 \right). \quad (12)$$

DP Principal Component Analysis

Lower bound for any data

DP-PCA: Statistically Optimal and Differentially Private PCA

Xiyang Liu * Weihao Kong † Prateek Jain ‡ Sewoong Oh §

Theorem 5.4 (Lower bound without Assumptions A.1–A.3). Let \mathcal{M}_ϵ be a set of DP estimators that map n i.i.d. samples to an estimate $\hat{v} \in \mathbb{R}^d$. A set of distributions satisfying Assumptions A.1–A.3 with $M = \tilde{O}(d + \sqrt{n\epsilon/d})$, $V = O(d)$ and $\gamma = O(1)$ is denoted by $\tilde{\mathcal{P}}$. For $d \geq 2$, there exists a universal constant $C > 0$ such that

$$\inf_{\hat{v} \in \mathcal{M}_\epsilon} \sup_{P \in \tilde{\mathcal{P}}} \mathbb{E}_{S \sim P^n} [\sin(\hat{v}(S), v_1)] \geq C \kappa \min \left(\sqrt{\frac{d \wedge \log((1 - e^{-\epsilon})/\delta)}{\epsilon n}}, 1 \right). \quad (13)$$

Theorem 5.5 (Lower bound, Gaussian distribution). Let \mathcal{M}_ϵ be a set of DP estimators that map n i.i.d. samples to an estimate $\hat{v} \in \mathbb{R}^d$. A set of Gaussian distributions with (λ_1, λ_2) as the first and second eigenvalues of the covariance matrix is denoted by $\mathcal{P}_{(\lambda_1, \lambda_2)}$. There exists a universal constant $C > 0$ such that

$$\inf_{\hat{v} \in \mathcal{M}_\epsilon} \sup_{P \in \mathcal{P}_{(\lambda_1, \lambda_2)}} \mathbb{E}_{S \sim P^n} [\sin(\hat{v}(S), v_1)] \geq C \min \left(\kappa \left(\sqrt{\frac{d}{n}} - \frac{d}{\epsilon n} \right) \sqrt{\frac{\lambda_2}{\lambda_1}}, 1 \right). \quad (12)$$

Lower bound for sub-gaussian

DP Principal Component Analysis

Lower bound for any data

DP-PCA: Statistically Optimal and Differentially Private PCA

Xiyang Liu * Weihao Kong † Prateek Jain ‡ Sewoong Oh §

Theorem 5.4 (Lower bound without Assumptions A.1–A.3). Let \mathcal{P} be a set of distributions that map n i.i.d. samples to an estimate $\hat{v} \in \mathbb{R}^d$. A set of distributions satisfying Assumptions A.1–A.3 with $\gamma = O(1)$ is denoted by $\tilde{\mathcal{P}}$. For $d \geq 2$, there exists a

$$C \kappa \min \left(\sqrt{\frac{d \wedge \log((1 - e^{-\varepsilon})/\delta)}{\varepsilon n}}, 1 \right). \quad (13)$$

- Similar results known about other problems including PCA.
- Key takeaway:
 - Privacy guaranteed for **ALL** datasets.
 - **When data quality is high, utility is better.**

A set of Gaussian distributions with (λ_1, λ_2) as the covariance matrix is denoted by $\mathcal{P}_{(\lambda_1, \lambda_2)}$. There exists a universal

$$\inf_{\hat{v} \in \mathcal{M}_\varepsilon} \sup_{P \in \mathcal{P}_{(\lambda_1, \lambda_2)}} \mathbb{E}_{S \sim P^n} [\sin(\hat{v}(S), v_1)] \geq C \min \left(\kappa \left(\sqrt{\frac{d}{n}} - \frac{d}{\varepsilon n} \right) \sqrt{\frac{\lambda_2}{\lambda_1}}, 1 \right). \quad (12)$$

Lower bound for sub-gaussian

DP SGD

Deep Learning with Differential Privacy

October 25, 2016

Martin Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow†
Kunal Talwar*

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martin Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow†
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

Take a random sample L_t with sampling probability L/N

Compute gradient

For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

Clip gradient

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martin Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability L/N

Compute gradient

For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

Clip gradient

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability L/N

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

Clip gradient

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability $\frac{1}{L/N}$

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability L/N

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**
 $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability L/N

(2) **Compute gradient**

For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

(5) **Descent**

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.
- As we saw earlier, the added noise scales with dimensionality of params

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow*
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability L/N

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**
 $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

(5) **Descent**
 $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.
- As we saw earlier, the added noise scales with dimensionality of params
- To avoid this, they conduct DP-PCA on data before doing DP-SGD.

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow†
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability $\frac{1}{L/N}$

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**
 $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

(5) **Descent**
 $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.
- As we saw earlier, the added noise scales with dimensionality of params
- To avoid this, they conduct DP-PCA on data before doing DP-SGD.

But,

Deep Learning with Differential Privacy

October 25, 2016

Martín Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow†
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability $\frac{1}{L/N}$

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**
 $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

(5) **Descent**
 $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.
- As we saw earlier, the added noise scales with dimensionality of params
- To avoid this, they conduct DP-PCA on data before doing DP-SGD.

But,

1. DP-PCA requires **additional time**

Deep Learning with Differential Privacy

October 25, 2016

Martin Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow†
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability $\frac{1}{L/N}$

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**
 $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

(5) **Descent**
 $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP SGD

- DP-SGD is the standard workhorse for DP Machine Learning algorithms.
- As we saw earlier, the added noise scales with dimensionality of params
- To avoid this, they conduct DP-PCA on data before doing DP-SGD.

But,

1. DP-PCA requires **additional time**
2. DP-PCA incurs **additional privacy cost**

Deep Learning with Differential Privacy

October 25, 2016

Martin Abadi*
H. Brendan McMahan*

Andy Chu*
Ilya Mironov*
Li Zhang*

Ian Goodfellow†
Kunal Talwar*

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

(1) Take a random sample L_t with sampling probability $\frac{1}{L/N}$

(2) **Compute gradient**
For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

(3) **Clip gradient**
 $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

(4) **Add noise**
 $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

(5) **Descent**
 $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

PILLAR

Leveraging intrinsic low dimensionality

2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk

PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk

PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk

PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space
- Maybe we need to find the right representation space

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk

PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space
- Maybe we need to find the right representation space

Idea 2: Use any **public unlabelled pre-training** data for representation learning.



PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk

PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space
- Maybe we need to find the right representation space

Idea 2: Use any **public unlabelled pre-training data** for representation learning.



2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

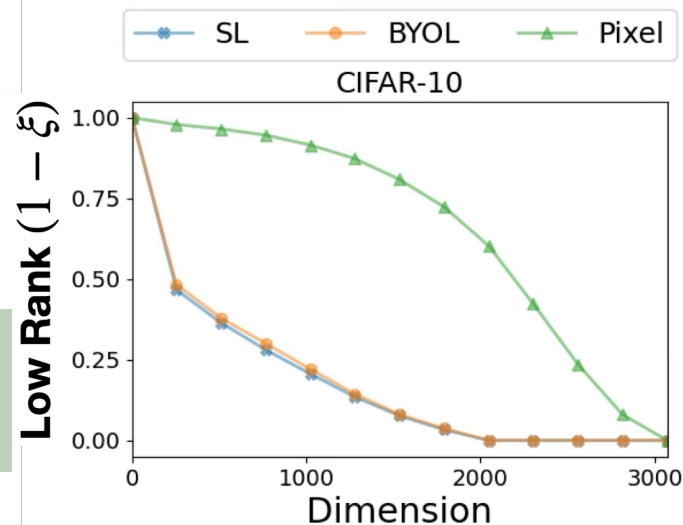
PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space
- Maybe we need to find the right representation space

Idea 2: Use any **public unlabelled pre-training data** for representation learning.

$1 - \xi$ = low rank reconstruction error



2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

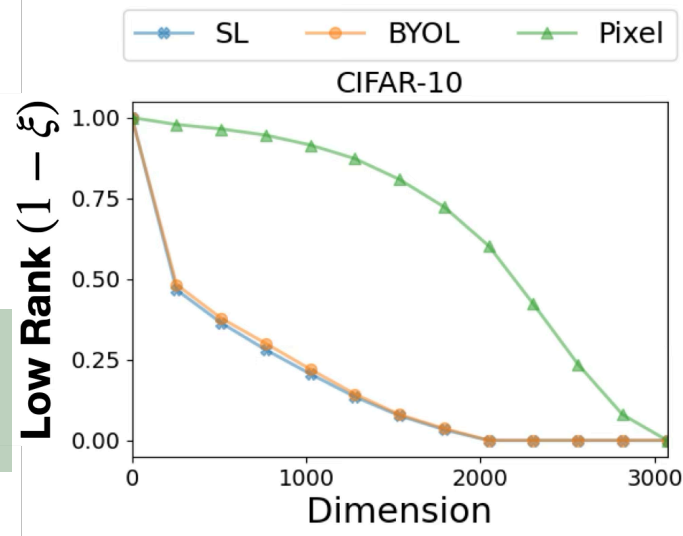
PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space
- Maybe we need to find the right representation space

Idea 2: Use any **public unlabelled pre-training data** for representation learning.

$1 - \xi$ = low rank reconstruction error



2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

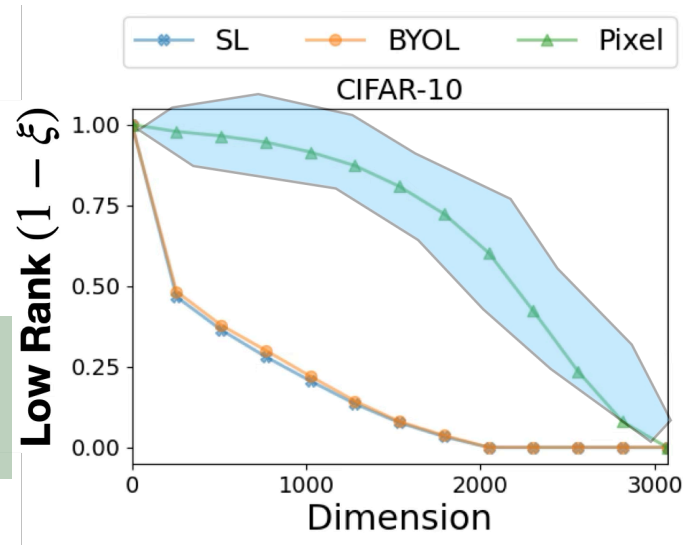
PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



PILLAR

Leveraging intrinsic low dimensionality

Idea 1: Use identically distributed **public unlabelled data** to find low rank subspace for projection



- But natural data is not inherently low rank in pixel space
- Maybe we need to find the right representation space

Idea 2: Use any **public unlabelled pre-training data** for representation learning.

$1 - \xi$ = low rank reconstruction error



2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

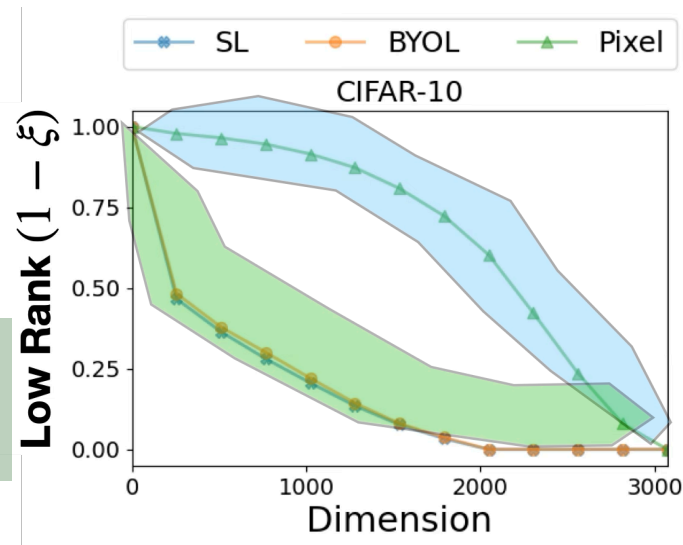
PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk

PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



Public
unlabelled
pre-training



PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



Public
unlabelled
pre-training

PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto*
University of Oxford
Oxford, England
francesco.pinto@eng.ox.ac.uk

Yaxi Hu*
Max Planck Institute for Intelligence Systems
Tübingen, Germany
yaxi.hu@tuebingen.mpg.de

Fanny Yang
ETH Zürich
Zürich, Switzerland
fan.yang@inf.ethz.ch

Amartya Sanyal
Max Planck Institute for Intelligence Systems
Tübingen, Germany
amsa@di.ku.dk



PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

Private labelled



2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto* University of Oxford Oxford, England francesco.pinto@eng.ox.ac.uk	Yaxi Hu* Max Planck Institute for Intelligence Systems Tübingen, Germany yaxi.hu@tuebingen.mpg.de
Fanny Yang ETH Zürich Zürich, Switzerland fanny.yang@inf.ethz.ch	Amartya Sanyal Max Planck Institute for Intelligence Systems Tübingen, Germany amsa@di.ku.dk



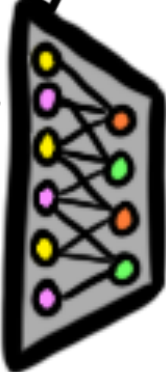
Public
unlabelled
pre-training



PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

Private labelled



Public
unlabelled
pre-training

2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto* University of Oxford Oxford, England francesco.pinto@eng.ox.ac.uk	Yaxi Hu* Max Planck Institute for Intelligence Systems Tübingen, Germany yaxi.hu@tuebingen.mpg.de
Fanny Yang ETH Zürich Zürich, Switzerland fan.yang@inf.ethz.ch	Amartya Sanyal Max Planck Institute for Intelligence Systems Tübingen, Germany amsa@di.ku.dk

PILLAR for Chest X-Ray Classification

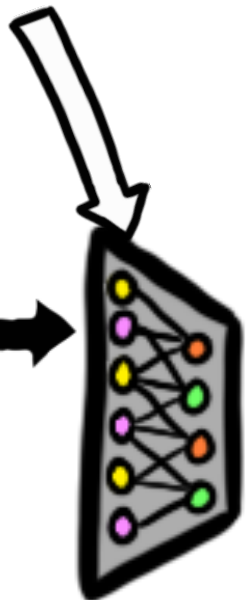
Leveraging intrinsic low dimensionality

Private labelled

Public unlabelled



Public unlabelled pre-training



2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto* University of Oxford Oxford, England francesco.pinto@eng.ox.ac.uk	Yaxi Hu* Max Planck Institute for Intelligence Systems Tübingen, Germany yaxi.hu@tuebingen.mpg.de
Fanny Yang ETH Zürich Zürich, Switzerland fan.yang@inf.ethz.ch	Amartya Sanyal Max Planck Institute for Intelligence Systems Tübingen, Germany amsa@di.ku.dk

PILLAR for Chest X-Ray Classification

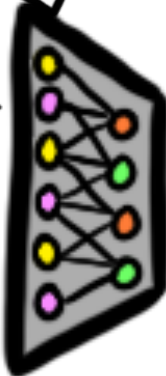
Leveraging intrinsic low dimensionality

Private labelled

Public unlabelled



Public unlabelled pre-training



PILLAR for Chest X-Ray Classification

Leveraging intrinsic low dimensionality

Private labelled

Public unlabelled

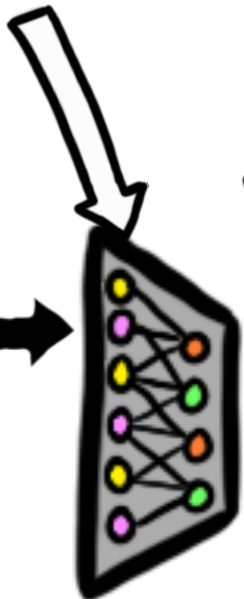
2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)

PILLAR: How to make semi-private learning more effective

Francesco Pinto* University of Oxford Oxford, England francesco.pinto@eng.ox.ac.uk	Yaxi Hu* Max Planck Institute for Intelligence Systems Tübingen, Germany yaxi.hu@tuebingen.mpg.de
Fanny Yang ETH Zürich Zürich, Switzerland fan.yang@inf.ethz.ch	Amartya Sanyal Max Planck Institute for Intelligence Systems Tübingen, Germany amsa@di.ku.dk



Public unlabelled pre-training



PILLAR

PILLAR for Chest X-Ray Classification

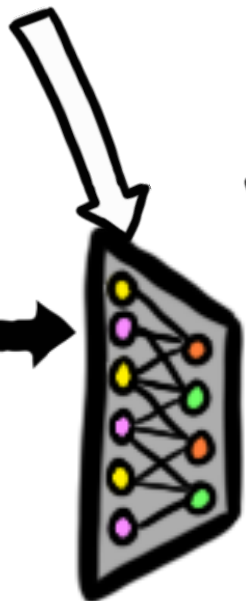
Leveraging intrinsic low dimensionality

Private labelled

Public unlabelled



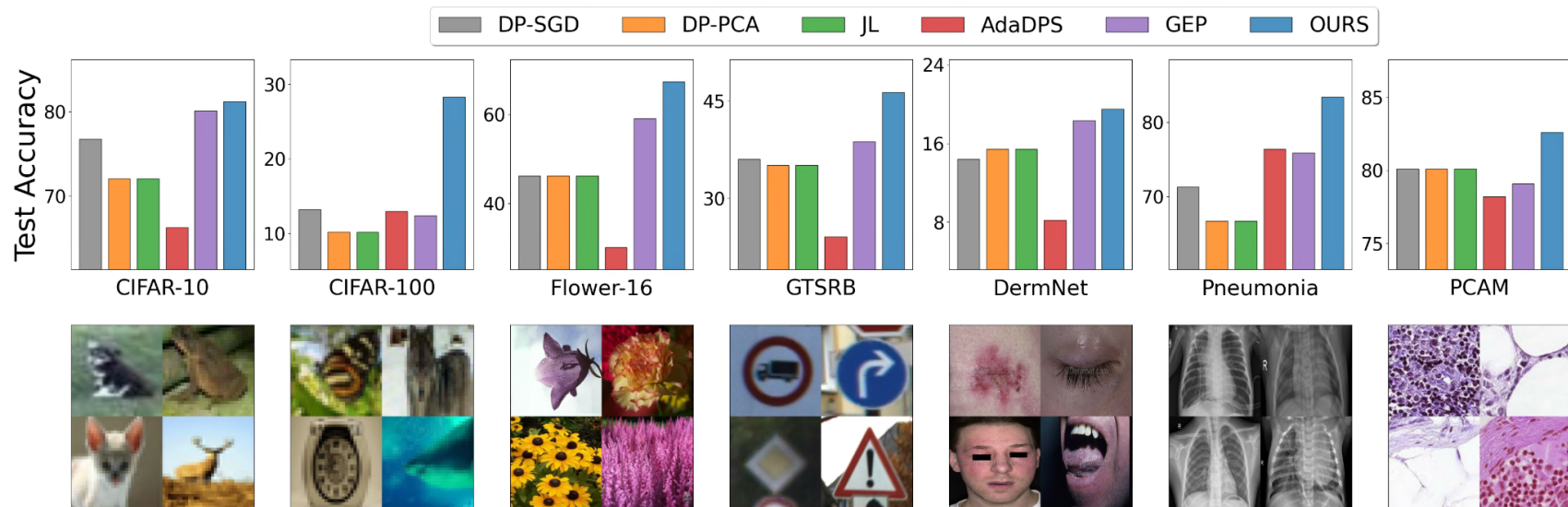
Public unlabelled pre-training



PILLAR



Other approaches to leverage unlabelled data



- GEP works in the gradient space
- AdaDPS use public data for gradient pre-conditioning

Next

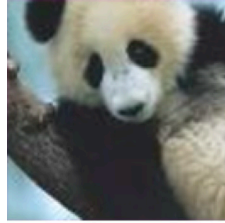


Robustness in Machine Learning

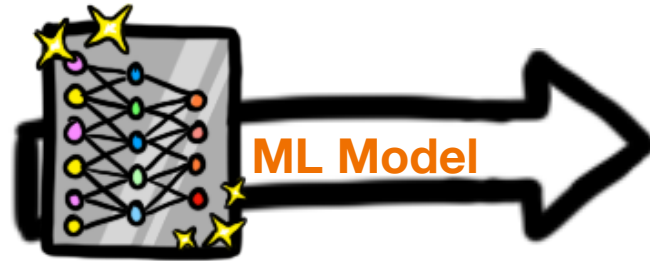
Adversarial Robustness in Machine Learning

Adversarial Robustness in Machine Learning

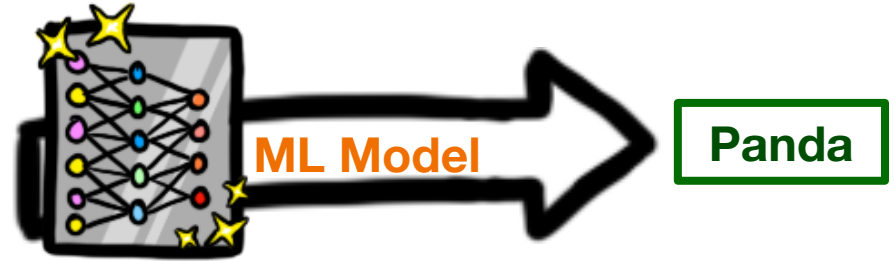
Adversarial Robustness in Machine Learning



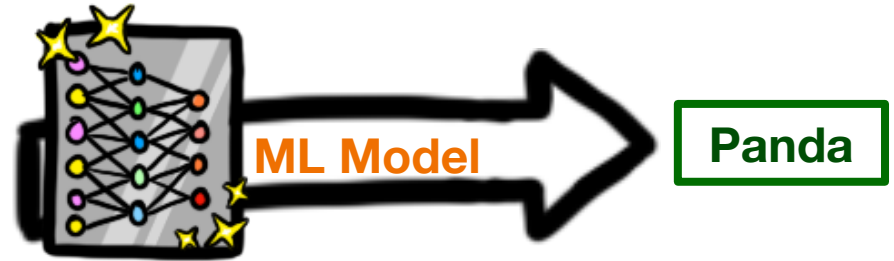
Adversarial Robustness in Machine Learning



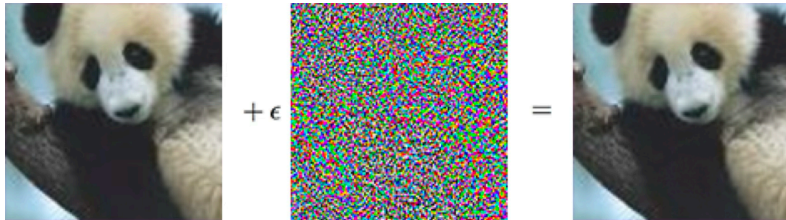
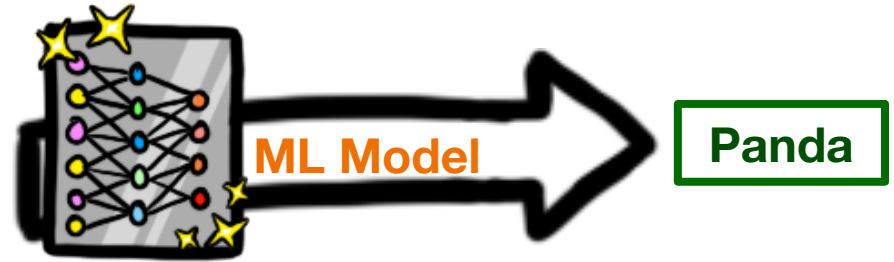
Adversarial Robustness in Machine Learning



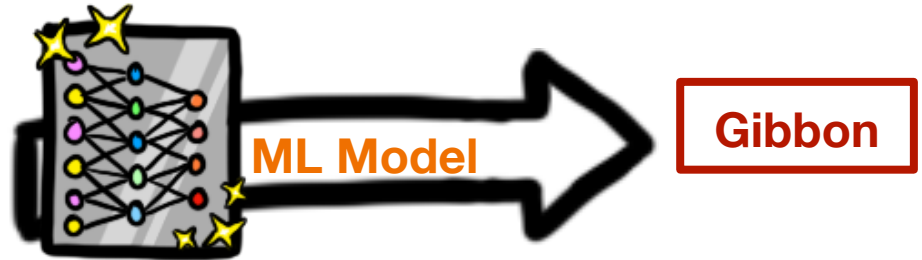
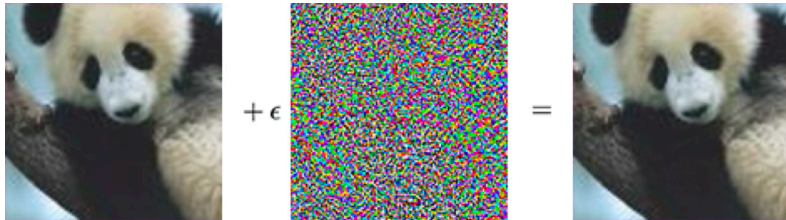
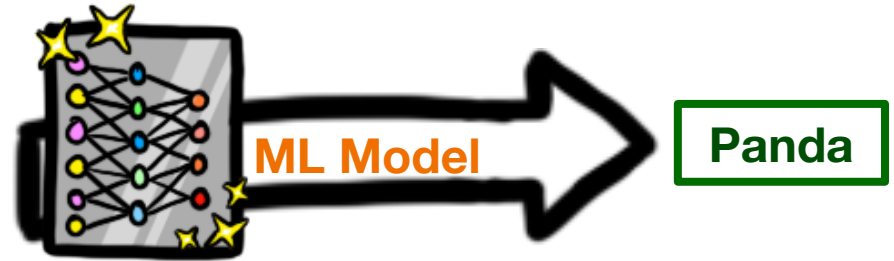
Adversarial Robustness in Machine Learning



Adversarial Robustness in Machine Learning



Adversarial Robustness in Machine Learning

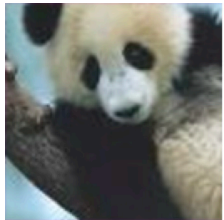


Adversarial Robustness in Machine Learning

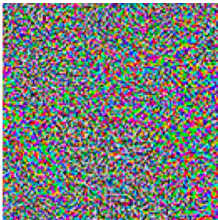


ML Model

Panda



+ ϵ



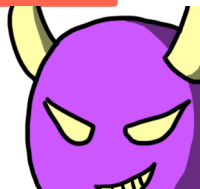
=



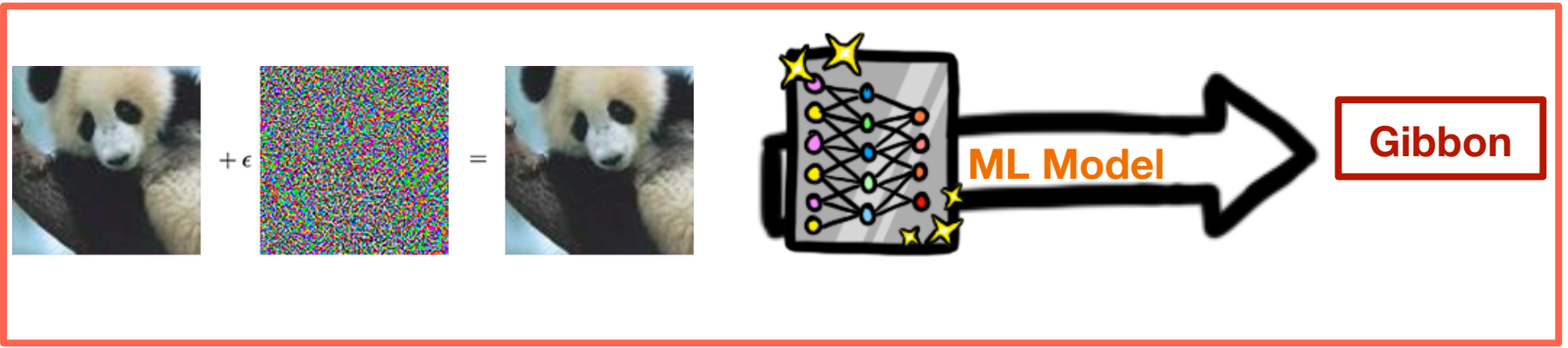
ML Model

Gibbon

Adversarial Example



Adversarial Robustness in Machine Learning



Adversarial Robustness in Machine Learning



For any distribution \mathcal{P} over $\mathbb{R}^d \times \{0, 1\}$ and any binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$

Adversarial Robustness in Machine Learning



For any distribution \mathcal{P} over $\mathbb{R}^d \times \{0, 1\}$ and any binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$ the γ -adversarial error is defined as

Adversarial Robustness in Machine Learning

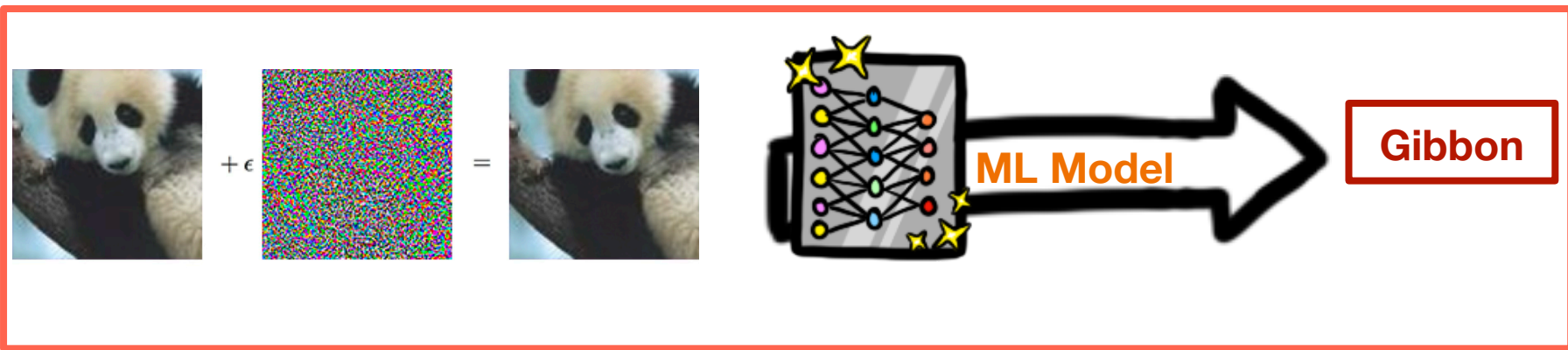


For any distribution \mathcal{P} over $\mathbb{R}^d \times \{0, 1\}$ and any binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$

the γ -adversarial error is defined as

$$\Pr_{(x,y) \sim \mathcal{P}} [\text{exists } z \in \mathcal{B}_\gamma(x) : f(z) \neq y]$$

Adversarial Robustness in Machine Learning



For any distribution \mathcal{P} over $\mathbb{R}^d \times \{0, 1\}$ and any binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$

the γ -adversarial error is defined as

$$\Pr_{(x,y) \sim \mathcal{P}} [\text{exists } z \in \mathcal{B}_\gamma(x) : f(z) \neq y]$$

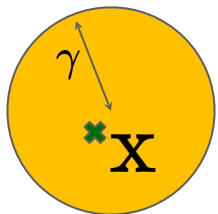
*x

Adversarial Robustness in Machine Learning



For any distribution \mathcal{P} over $\mathbb{R}^d \times \{0, 1\}$ and any binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$

the γ -adversarial error is defined as



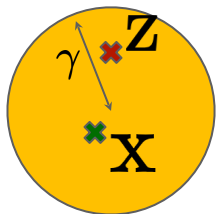
$$\Pr_{(\mathbf{x}, y) \sim \mathcal{P}} [\text{exists } \mathbf{z} \in \mathcal{B}_\gamma(\mathbf{x}) : f(\mathbf{z}) \neq y]$$

Adversarial Robustness in Machine Learning



For any distribution \mathcal{P} over $\mathbb{R}^d \times \{0, 1\}$ and any binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$

the γ -adversarial error is defined as



$$\Pr_{(x,y) \sim \mathcal{P}} [\text{exists } z \in \mathcal{B}_\gamma(x) : f(z) \neq y]$$

Label Noise is ubiquitously interpolated

UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minsok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Label Noise is ubiquitously interpolated

- Trained long enough, NNs fit label noise

UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

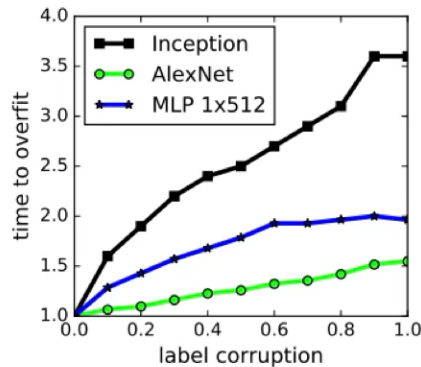
Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minsok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Label Noise is ubiquitously interpolated

- Trained long enough, NNs fit label noise



UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

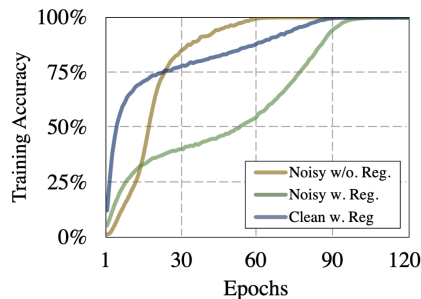
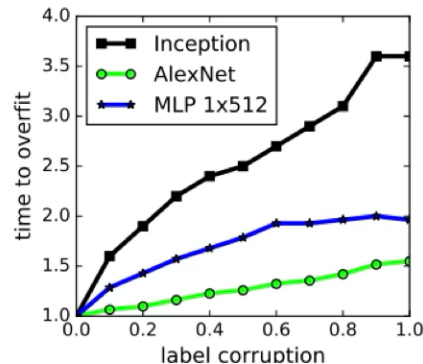
Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minsok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Label Noise is ubiquitously interpolated

- Trained long enough, NNs fit label noise



UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang^{*}
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht[†]
University of California, Berkeley
brecht@berkeley.edu

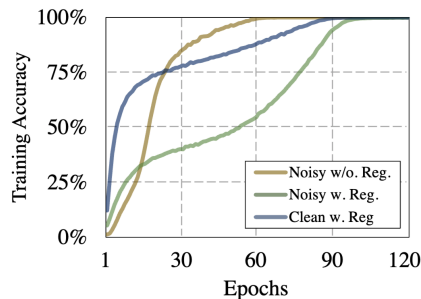
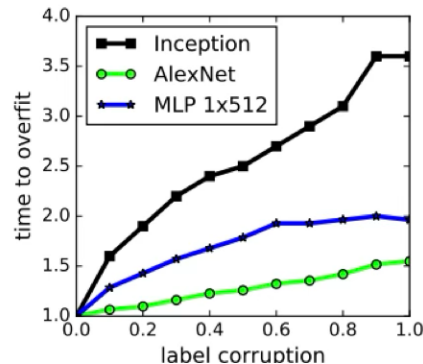
Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Label Noise is ubiquitously interpolated

- Trained long enough, NNs fit label noise
- Does not always hurt Test Accuracy - Benign Overfitting



UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

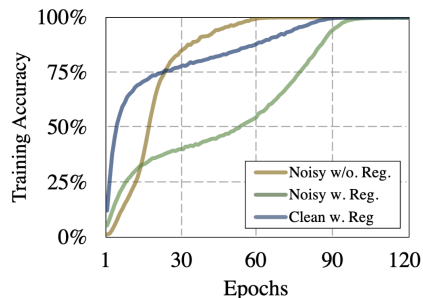
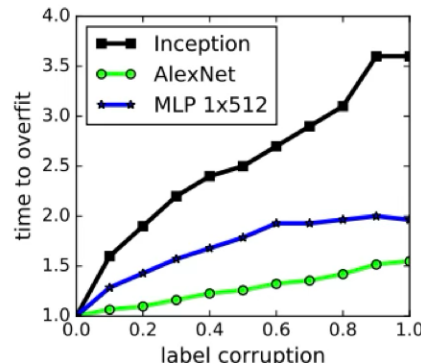
Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Label Noise is ubiquitously interpolated

- Trained long enough, NNs fit label noise
- Does not always hurt Test Accuracy - Benign Overfitting
- Define a model with 100% training acc: **Interpolator**



UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

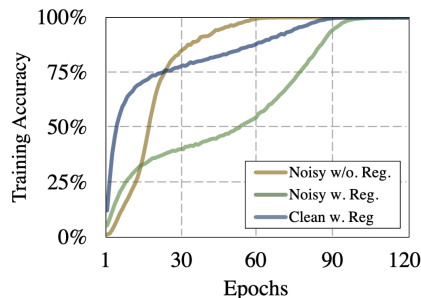
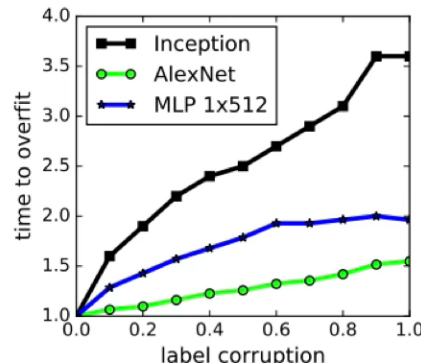
Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Label Noise is ubiquitously interpolated

- Trained long enough, NNs fit label noise
- Does not always hurt Test Accuracy - Benign Overfitting
- Define a model with 100% training acc: **Interpolator**



UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

Oriol Vinyals
Google DeepMind
vinyals@google.com

Learning from Noisy Labels with Deep Neural Networks: A Survey

Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, Jae-Gil Lee

Question: What about **Robust Accuracy** ?

Lower bound on Adversarial error

HOW BENIGN IS BENIGN OVERFITTING?

Amartya Sanyal
Department of Computer Science,
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
amartya.sanyal@cs.ox.ac.uk

Varun Kanade
Department of Computer Science
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
varunk@cs.ox.ac.uk

Puneet K. Dokania
Department of Engineering Science
University of Oxford, Oxford, UK
Five AI Limited
puneet@robots.ox.ac.uk

Philip H.S. Torr
Department of Engineering Science
University of Oxford, Oxford, UK
phst@robots.ox.ac.uk

A LAW OF ADVERSARIAL RISK, INTERPOLATION, AND LABEL NOISE

Daniel Paleka *
ETH Zurich
daniel.paleka@inf.ethz.ch

Amartya Sanyal *
ETH AI Center, ETH Zurich
amartya.sanyal@ai.ethz.ch

Lower bound on Adversarial error

Let

- μ be any distribution on \mathbb{R}^d ,
- $\eta \in (0, 1)$ be the uniform label noise rate,
- $\mathcal{C} \subset \mathbb{R}^d$ be any region, and
- $N(\mathcal{C}, \epsilon, \|\cdot\|)$ is the covering number of \mathcal{C}

HOW BENIGN IS BENIGN OVERFITTING?

Amartya Sanyal
Department of Computer Science,
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
amartya.sanyal@cs.ox.ac.uk

Varun Kanade
Department of Computer Science
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
varunk@cs.ox.ac.uk

Puneet K. Dokania
Department of Engineering Science
University of Oxford, Oxford, UK
Five AI Limited
puneet@robots.ox.ac.uk

Philip H.S. Torr
Department of Engineering Science
University of Oxford, Oxford, UK
phst@robots.ox.ac.uk

A LAW OF ADVERSARIAL RISK, INTERPOLATION, AND LABEL NOISE

Daniel Paleka *
ETH Zurich
daniel.paleka@inf.ethz.ch

Amartya Sanyal *
ETH AI Center, ETH Zurich
amartya.sanyal@ai.ethz.ch

Lower bound on Adversarial error

Let

- μ be any distribution on \mathbb{R}^d ,
- $\eta \in (0, 1)$ be the uniform label noise rate,
- $\mathcal{C} \subset \mathbb{R}^d$ be any region, and
- $N(\mathcal{C}, \epsilon, \|\cdot\|)$ is the covering number of \mathcal{C}

HOW BENIGN IS BENIGN OVERFITTING?

Amartya Sanyal
Department of Computer Science,
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
amartya.sanyal@cs.ox.ac.uk

Varun Kanade
Department of Computer Science
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
varunk@cs.ox.ac.uk

Puneet K. Dokania
Department of Engineering Science
University of Oxford, Oxford, UK
Five AI Limited
puneet@robots.ox.ac.uk

Philip H.S. Torr
Department of Engineering Science
University of Oxford, Oxford, UK
phst@robots.ox.ac.uk

A LAW OF ADVERSARIAL RISK, INTERPOLATION, AND LABEL NOISE

Daniel Paleka *
ETH Zurich
daniel.paleka@inf.ethz.ch

Amartya Sanyal *
ETH AI Center, ETH Zurich
amartya.sanyal@ai.ethz.ch

Theorem If the noisy dataset size $m = \Omega\left(\frac{N(\mathcal{C}, \epsilon, \|\cdot\|)}{\mu(\mathcal{C})\eta}\right)$, for all interpolators h

Lower bound on Adversarial error

Let

- μ be any distribution on \mathbb{R}^d ,
- $\eta \in (0, 1)$ be the uniform label noise rate,
- $\mathcal{C} \subset \mathbb{R}^d$ be any region, and
- $N(\mathcal{C}, \epsilon, \|\cdot\|)$ is the covering number of \mathcal{C}

HOW BENIGN IS BENIGN OVERFITTING?

Amartya Sanyal
Department of Computer Science,
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
amartya.sanyal@cs.ox.ac.uk

Varun Kanade
Department of Computer Science
University of Oxford,
Oxford, UK
The Alan Turing Institute, London, UK
varunk@cs.ox.ac.uk

Puneet K. Dokania
Department of Engineering Science
University of Oxford, Oxford, UK
Five AI Limited
puneet@robots.ox.ac.uk

Philip H.S. Torr
Department of Engineering Science
University of Oxford, Oxford, UK
phst@robots.ox.ac.uk

A LAW OF ADVERSARIAL RISK, INTERPOLATION, AND LABEL NOISE

Daniel Paleka *
ETH Zurich
daniel.paleka@inf.ethz.ch

Amartya Sanyal *
ETH AI Center, ETH Zurich
amartya.sanyal@ai.ethz.ch

Theorem If the noisy dataset size $m = \Omega\left(\frac{N(\mathcal{C}, \epsilon, \|\cdot\|)}{\mu(\mathcal{C})\eta}\right)$, for all interpolators h

$$\text{Adv. Error}_\epsilon(h) \geq \mu(\mathcal{C})$$

Lower bound on Adversarial error

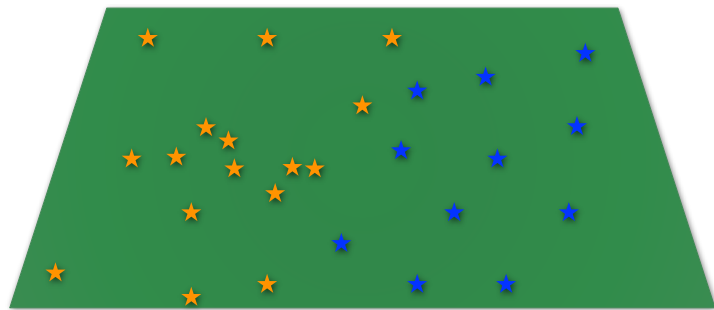
Lower bound on Adversarial error



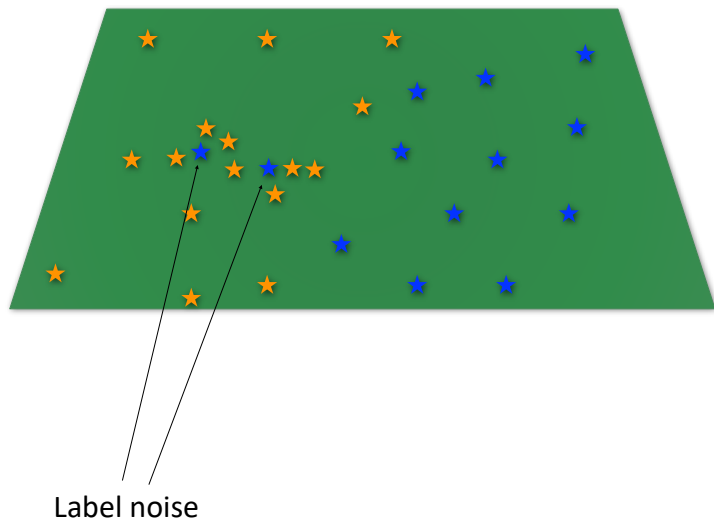
Lower bound on Adversarial error



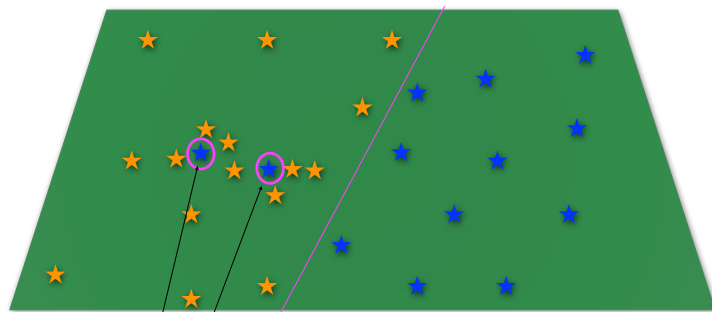
Lower bound on Adversarial error



Lower bound on Adversarial error



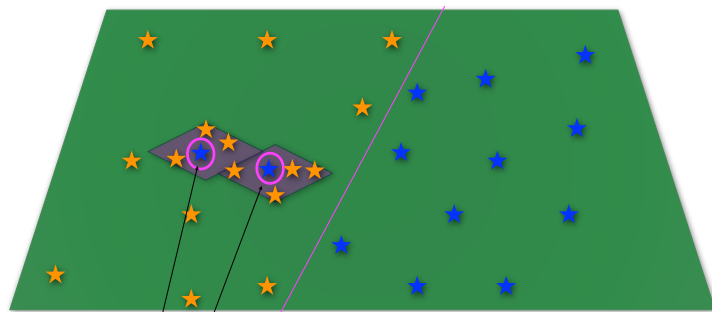
Lower bound on Adversarial error



Label noise

h is an interpolator e.g. Random Forest, 1-NN, NNs

Lower bound on Adversarial error



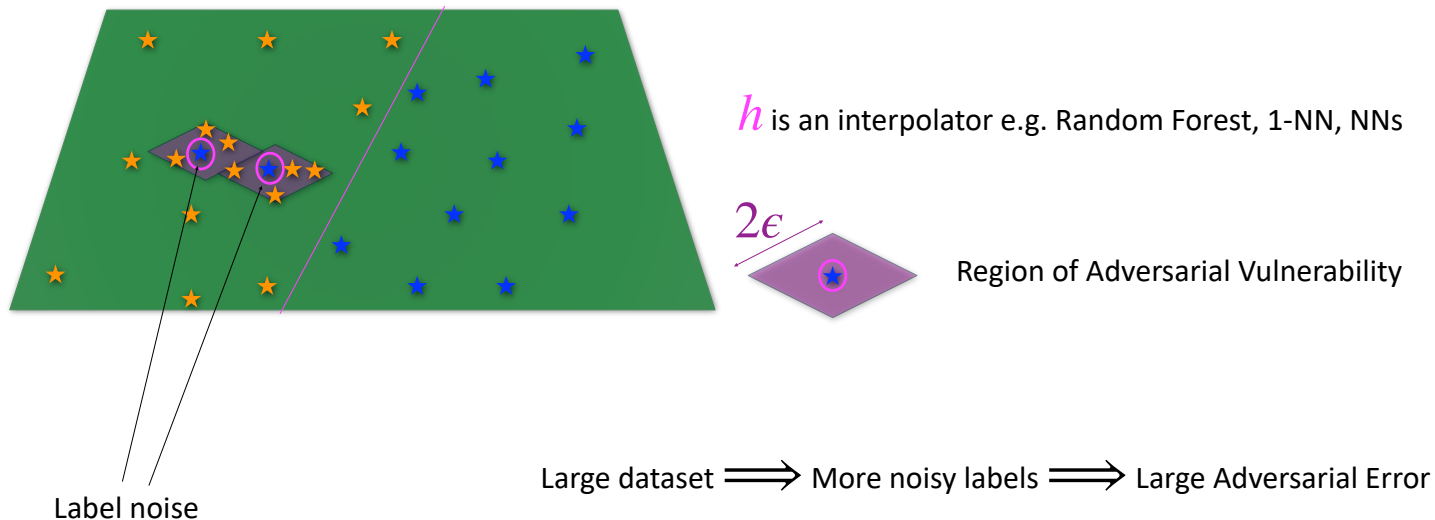
Label noise

h is an interpolator e.g. Random Forest, 1-NN, NNs



Region of Adversarial Vulnerability

Lower bound on Adversarial error



Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Madry*

MIT

madry@mit.edu

Aleksandar Makelov*

MIT

amakelov@mit.edu

Ludwig Schmidt*

MIT

ludwigs@mit.edu

Dimitris Tsipras*

MIT

tsipras@mit.edu

Adrian Vladu*

MIT

avladu@mit.edu

Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Mądry* MIT madr@mit.edu	Aleksandar Makelov* MIT amakelov@mit.edu	Ludwig Schmidt* MIT ludwigs@mit.edu
Dimitris Tsipras* MIT tsipras@mit.edu	Adrian Vladu* MIT avladu@mit.edu	

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.

Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Mądry*
MIT
mardy@mit.edu

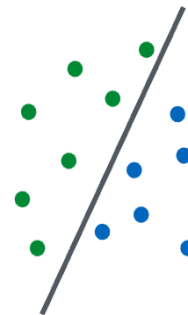
Aleksandar Makelov*
MIT
amakelov@mit.edu

Ludwig Schmidt*
MIT
ludwigs@mit.edu

Dimitris Tsipras*
MIT
tsipras@mit.edu

Adrian Vladu*
MIT
avladu@mit.edu

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.

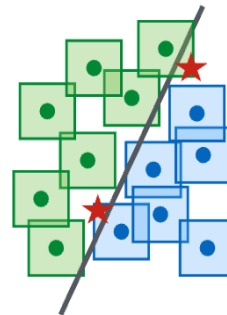


Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Mądry* MIT madr@mit.edu	Aleksandar Makelov* MIT amakelov@mit.edu	Ludwig Schmidt* MIT ludwigs@mit.edu
Dimitris Tsipras* MIT tsipras@mit.edu	Adrian Vladu* MIT avladu@mit.edu	

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.

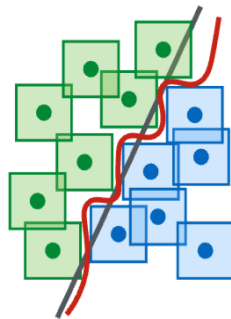


Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Mądry* MIT madr@mit.edu	Aleksandar Makelov* MIT amakelov@mit.edu	Ludwig Schmidt* MIT ludwigs@mit.edu
Dimitris Tsipras* MIT tsipras@mit.edu	Adrian Vladu* MIT avladu@mit.edu	

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.



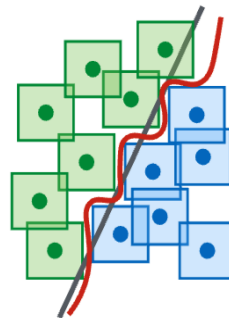
Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Mądry* MIT madr@mit.edu	Aleksandar Makelov* MIT amakelov@mit.edu	Ludwig Schmidt* MIT ludwigs@mit.edu
Dimitris Tsipras* MIT tsipras@mit.edu	Adrian Vladu* MIT avladu@mit.edu	

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.

Naturally, complex models can fit the augmented data better.



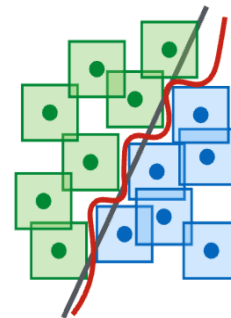
Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Mądry* MIT madr@mit.edu	Aleksandar Makelov* MIT amakelov@mit.edu	Ludwig Schmidt* MIT ludwigs@mit.edu
Dimitris Tsipras* MIT tsipras@mit.edu	Adrian Vladu* MIT avladu@mit.edu	

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.

Naturally, complex models can fit the augmented data better.



Robust overfitting is when train robust error decreases but test robust error increases.

Overfitting in adversarially robust deep learning

Leslie Rice^{*1} Eric Wong^{*2} J. Zico Kolter¹

Adversarial training

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Madry*
MIT
madry@mit.edu

Aleksandar Makelov*
MIT
amakelov@mit.edu

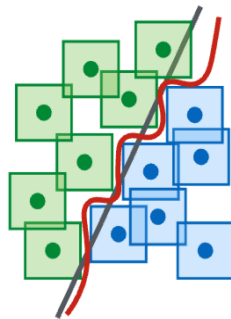
Ludwig Schmidt*
MIT
ludwigs@mit.edu

Dimitris Tsipras*
MIT
tsipras@mit.edu

Adrian Vladu*
MIT
avladu@mit.edu

Adversarial Training replaces (or augments) clean data with corresponding adversarial examples during SGD.

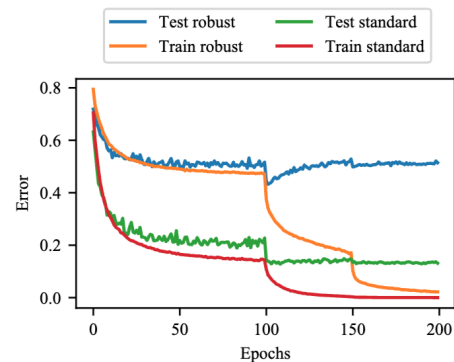
Naturally, complex models can fit the augmented data better.



Robust overfitting is when train robust error decreases but test robust error increases.

Overfitting in adversarially robust deep learning

Leslie Rice^{*1} Eric Wong^{*2} J. Zico Kolter¹



Robust Overfitting and label noise

**Label Noise in Adversarial Training: A Novel
Perspective to Study Robust Overfitting**

Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

Liyuan Liu
Microsoft Research
lucliu@microsoft.com

Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu

Robust Overfitting and label noise

- One of explanations given for Robust overfitting is that adversarial training implicitly adds label noise.

Label Noise in Adversarial Training: A Novel Perspective to Study Robust Overfitting

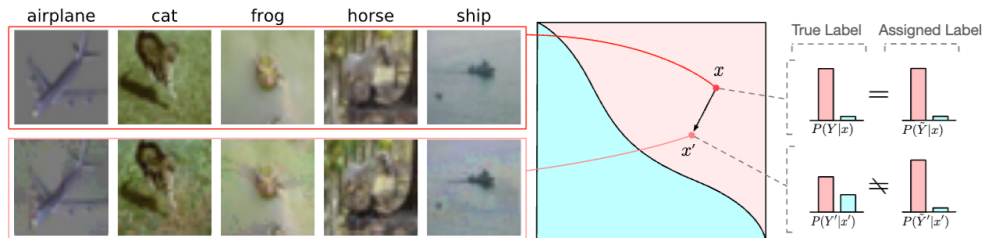
Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

Liyuan Liu
Microsoft Research
lucliu@microsoft.com

Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu

Robust Overfitting and label noise

- One of explanations given for Robust overfitting is that adversarial training implicitly adds label noise.



Label Noise in Adversarial Training: A Novel Perspective to Study Robust Overfitting

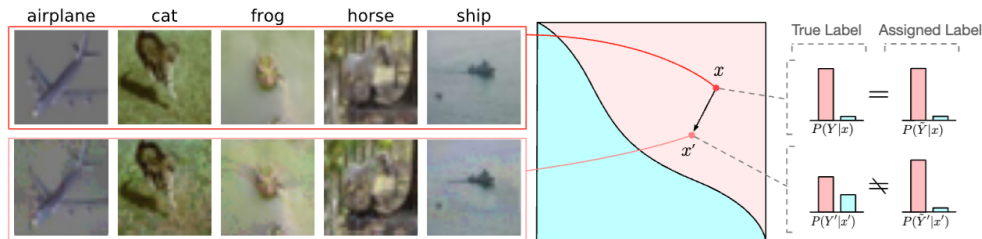
Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

Liyuan Liu
Microsoft Research
liuliu@microsoft.com

Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu

Robust Overfitting and label noise

- One of explanations given for Robust overfitting is that adversarial training implicitly adds label noise.



Label Noise in Adversarial Training: A Novel Perspective to Study Robust Overfitting

Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

Liyuan Liu
Microsoft Research
liuliu@microsoft.com

Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu

Data Quality Matters For Adversarial Training: An Empirical Study

Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

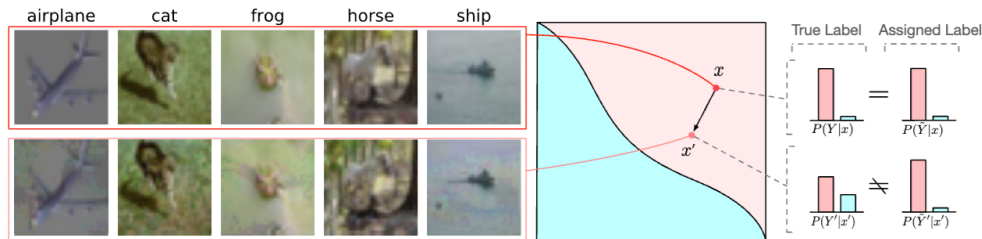
Liyuan Liu
University of Illinois at Urbana-Champaign
112@illinois.edu

Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu

- Simply using “good” examples that are far from the decision boundary alleviates parts of the issue

Robust Overfitting and label noise

- One of explanations given for Robust overfitting is that adversarial training implicitly adds label noise.



- Simply using “good” examples that are far from the decision boundary alleviates parts of the issue
- Larger perturbation radius causes more overfitting

Label Noise in Adversarial Training: A Novel Perspective to Study Robust Overfitting

Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

Liyuan Liu
Microsoft Research
liuliu@microsoft.com

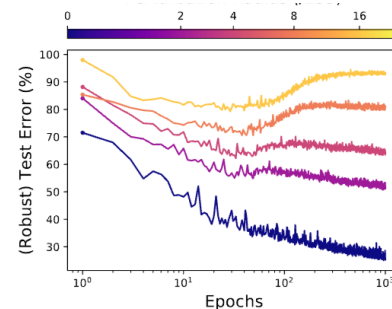
Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu

Data Quality Matters For Adversarial Training: An Empirical Study

Chengyu Dong
University of California, San Diego
cdong@eng.ucsd.edu

Liyuan Liu
University of Illinois at Urbana-Champaign
ll2@illinois.edu

Jingbo Shang
University of California, San Diego
jshang@eng.ucsd.edu



Improving Robust generalisation

Adversarially Robust Generalization Requires More Data

Ludwig Schmidt
MIT

Shibani Santurkar
MIT

Dimitris Tsipras
MIT

Kunal Talwar
Google Brain

Aleksander Madry
MIT

Interpolation can hurt robust generalization even when
there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
Reinhard Heckel^{2,3} and Fanny Yang¹

¹*ETH Zurich*, ²*Rice University*, ³*Technical University of Munich*

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

Adversarially Robust Generalization Requires More Data

Ludwig Schmidt
MIT

Shibani Santurkar
MIT

Dimitris Tsipras
MIT

Kunal Talwar
Google Brain

Aleksander Mądry
MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

- Clearly, more data helps to avoid robust overfitting

Adversarially Robust Generalization Requires More Data

Ludwig Schmidt
MIT

Shibani Santurkar
MIT

Dimitris Tsipras
MIT

Kunal Talwar
Google Brain

Aleksander Mądry
MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

- Clearly, more data helps to avoid robust overfitting
- Regularisation and early-stopping also helps.

Adversarially Robust Generalization Requires More Data

Ludwig Schmidt
MIT

Shibani Santurkar
MIT

Dimitris Tsipras
MIT

Kunal Talwar
Google Brain

Aleksander Mądry
MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

- Clearly, more data helps to avoid robust overfitting
- Regularisation and early-stopping also helps.

Adversarially Robust Generalization Requires More Data

Ludwig Schmidt
MIT

Shibani Santurkar
MIT

Dimitris Tsipras
MIT

Kunal Talwar
Google Brain

Aleksander Mądry
MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

- Clearly, more data helps to avoid robust overfitting
- Regularisation and early-stopping also helps.

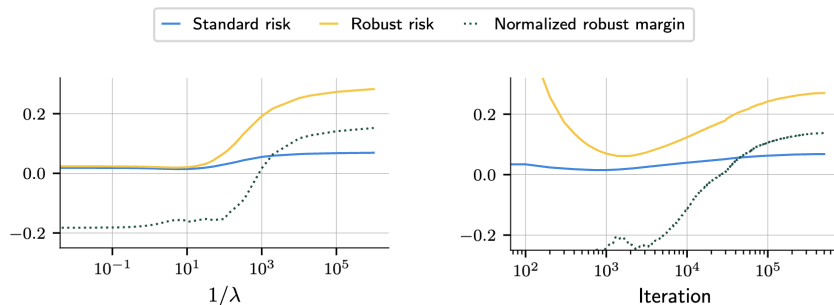
Adversarially Robust Generalization Requires More Data

Ludwig Schmidt MIT
Shibani Santurkar MIT
Dimitris Tsipras MIT
Kunal Talwar Google Brain
Aleksander Mądry MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich



(a) Benefit of ridge regularization

(b) Benefit of early stopping

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

- Clearly, more data helps to avoid robust overfitting
- Regularisation and early-stopping also helps.

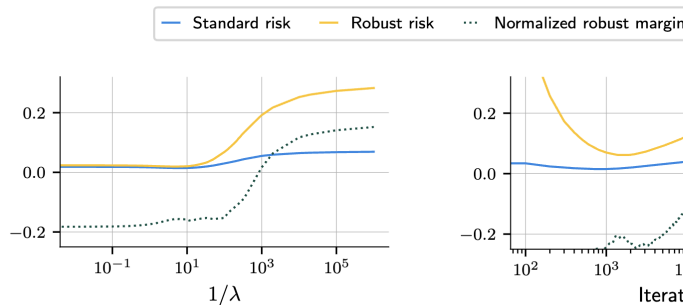
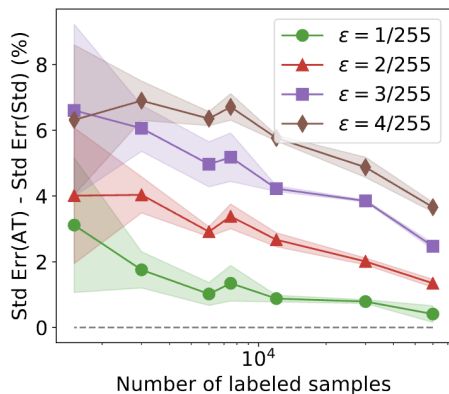
Adversarially Robust Generalization Requires More Data

Ludwig Schmidt MIT
 Kunal Talwar Google Brain
 Shibani Santurkar MIT
 Aleksander Mądry MIT
 Dimitris Tsipras MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
 Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich



(a) Benefit of ridge regularization

(b) Benefit of early stopping

Improving Robust generalisation

Exists simple distribution in d dim where robust generalisation requires \sqrt{d} times more data.

- Clearly, more data helps to avoid robust overfitting
- Regularisation and early-stopping also helps.

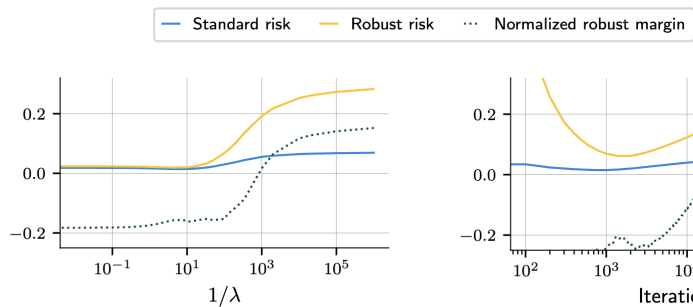
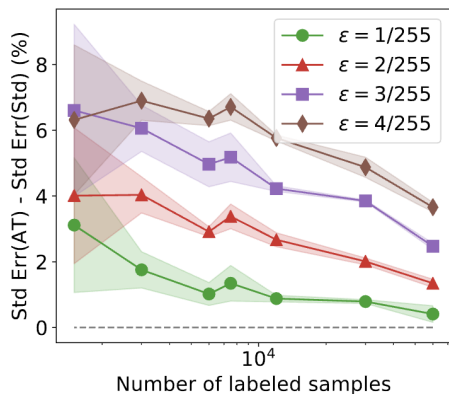
Adversarially Robust Generalization Requires More Data

Ludwig Schmidt MIT
 Kunal Talwar Google Brain
 Shibani Santurkar MIT
 Aleksander Mądry MIT
 Dimitris Tsipras MIT

Interpolation can hurt robust generalization even when there is no noise

Konstantin Donhauser¹, Alexandru Tifrea¹, Michael Aerni¹
 Reinhard Heckel^{2,3} and Fanny Yang¹

¹ETH Zurich, ²Rice University, ³Technical University of Munich



(a) Benefit of ridge regularization

(b) Benefit of early stopping

What about unlabelled data ?

Improving Robust generalisation

With unlabelled data

ADVERSARIALLY ROBUST GENERALIZATION JUST
REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹

¹Peking University ²Carnegie Mellon University ³Cornell University

⁴Huazhong University of Science and Technology

{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon* Aditi Raghunathan* Ludwig Schmidt
Stanford University Stanford University UC Berkeley
yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
Stanford University Stanford University
pliang@cs.stanford.edu jduchi@stanford.edu

Improving Robust generalisation

With unlabelled data

Observation: Robust error can be decomposed into

ADVERSARIALLY ROBUST GENERALIZATION JUST
REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹

¹Peking University ²Carnegie Mellon University ³Cornell University

⁴Huazhong University of Science and Technology

{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn

cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon* Aditi Raghunathan* Ludwig Schmidt
Stanford University Stanford University UC Berkeley
yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
Stanford University Stanford University
pliang@cs.stanford.edu jduchi@stanford.edu

Improving Robust generalisation

With unlabelled data

Observation: Robust error can be decomposed into

1. **Stability error:** Whether prediction is stable in a ball around data from the test distribution

ADVERSARIALLY ROBUST GENERALIZATION JUST
REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹

¹Peking University ²Carnegie Mellon University ³Cornell University

⁴Huazhong University of Science and Technology

{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn

cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon*
Stanford University
yairc@stanford.edu

Aditi Raghunathan*
Stanford University
aditir@stanford.edu

Ludwig Schmidt
UC Berkeley
ludwig@berkeley.edu

Percy Liang
Stanford University
pliang@cs.stanford.edu

John C. Duchi
Stanford University
jduchi@stanford.edu

Improving Robust generalisation

With unlabelled data

Observation: Robust error can be decomposed into

1. **Stability error:** Whether prediction is stable in a ball around data from the test distribution
2. **Classification accuracy:** Whether classification in the original data distribution is accurate

ADVERSARIALLY ROBUST GENERALIZATION JUST
REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹

¹Peking University ²Carnegie Mellon University ³Cornell University

⁴Huazhong University of Science and Technology

{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn

cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon*
Stanford University
yairc@stanford.edu

Aditi Raghunathan*
Stanford University
aditir@stanford.edu

Ludwig Schmidt
UC Berkeley
ludwig@berkeley.edu

Percy Liang
Stanford University
pliang@cs.stanford.edu

John C. Duchi
Stanford University
jduchi@stanford.edu

Improving Robust generalisation

With unlabelled data

Observation: Robust error can be decomposed into

- 1. Stability error:** Whether prediction is stable in a ball around data from the test distribution
- 2. Classification accuracy:** Whether classification in the original data distribution is accurate

Classical use of unlabelled data improves **2. Classification accuracy.**

ADVERSARIALLY ROBUST GENERALIZATION JUST
REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹
¹Peking University ²Carnegie Mellon University ³Cornell University
⁴Huazhong University of Science and Technology
{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon* Aditi Raghunathan* Ludwig Schmidt
Stanford University Stanford University UC Berkeley
yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
Stanford University Stanford University
pliang@cs.stanford.edu jduchi@stanford.edu

Improving Robust generalisation

With unlabelled data

Observation: Robust error can be decomposed into

- 1. Stability error:** Whether prediction is stable in a ball around data from the test distribution
- 2. Classification accuracy:** Whether classification in the original data distribution is accurate

ADVERSARIALLY ROBUST GENERALIZATION JUST REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹
¹Peking University ²Carnegie Mellon University ³Cornell University
⁴Huazhong University of Science and Technology
{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon* Aditi Raghunathan* Ludwig Schmidt
Stanford University Stanford University UC Berkeley
yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
Stanford University Stanford University
pliang@cs.stanford.edu jduchi@stanford.edu

Classical use of unlabelled data improves **2. Classification accuracy.**

To improve robustness, use unlabelled data to improve **1. Stability error.**

Improving Robust generalisation

With unlabelled data

Observation: Robust error can be decomposed into

1. **Stability error:** Whether prediction is stable in a ball around data from the test distribution
2. **Classification accuracy:** Whether classification in the original data distribution is accurate

ADVERSARIALLY ROBUST GENERALIZATION JUST REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹
¹Peking University ²Carnegie Mellon University ³Cornell University
⁴Huazhong University of Science and Technology
{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon* Aditi Raghunathan* Ludwig Schmidt
Stanford University Stanford University UC Berkeley
yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
Stanford University Stanford University
pliang@cs.stanford.edu jduchi@stanford.edu

Classical use of unlabelled data improves **2. Classification accuracy.**

To improve robustness, use unlabelled data to improve **1. Stability error.**

Recipe: Use adversarial training on pseudo-labels on the unlabelled data

Improving Robust generalisation

With unlabelled data

Improving Robust generalisation

With unlabelled data

Understanding and Mitigating the Tradeoff Between Robustness and Accuracy

Aditi Raghunathan^{*1} Sang Michael Xie^{*1} Fanny Yang² John C. Duchi¹ Percy Liang¹

Method	Robust Test Acc.	Standard Test Acc.	
Standard Training	0.8%	95.2%	} Vanilla Supervised
PG-AT (Madry et al., 2018)	45.8%	87.3%	
TRADES (Zhang et al., 2019)	55.4%	84.0%	
Standard Self-Training	0.3%	96.4%	} Semisupervised with same unlabeled data
Robust Consistency Training (Carmon et al., 2019)	56.5%	83.2%	
RST + PG-AT (this paper)	58.5%	91.8%	
RST + TRADES (this paper) (Carmon et al., 2019)	63.1%	89.7%	

Improving Robust generalisation

With unlabelled data

Understanding and Mitigating the Tradeoff Between Robustness and Accuracy

Aditi Raghunathan^{*1} Sang Michael Xie^{*1} Fanny Yang² John C. Duchi¹ Percy Liang¹

Method	Robust Test Acc.	Standard Test Acc.	
Standard Training	0.8%	95.2%	} Vanilla Supervised
PG-AT (Madry et al., 2018)	45.8%	87.3%	
TRADES (Zhang et al., 2019)	55.4%	84.0%	
Standard Self-Training	0.3%	96.4%	} Semisupervised with same unlabeled data
Robust Consistency Training (Carmon et al., 2019)	56.5%	83.2%	
RST + PG-AT (this paper)	58.5%	91.8%	
RST + TRADES (this paper) (Carmon et al., 2019)	63.1%	89.7%	

ADVERSARIALLY ROBUST GENERALIZATION JUST REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan²,
Kun He⁴, John E. Hopcroft³ & Liwei Wang¹
¹Peking University ²Carnegie Mellon University ³Cornell University
⁴Huazhong University of Science and Technology
{zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Improving Robust generalisation

With unlabelled data

Understanding and Mitigating the Tradeoff Between Robustness and Accuracy

Aaditi Raghunathan^{*1} Sang Michael Xie^{*1} Fanny Yang² John C. Duchi¹ Percy Liang¹

Method	Robust Test Acc.	Standard Test Acc.	
Standard Training	0.8%	95.2%	} Vanilla Supervised
PG-AT (Madry et al., 2018)	45.8%	87.3%	
TRADES (Zhang et al., 2019)	55.4%	84.0%	
Standard Self-Training	0.3%	96.4%	} Semisupervised with same unlabeled data
Robust Consistency Training (Carmon et al., 2019)	56.5%	83.2%	
RST + PG-AT (this paper)	58.5%	91.8%	
RST + TRADES (this paper) (Carmon et al., 2019)	63.1%	89.7%	

ADVERSARIALLY ROBUST GENERALIZATION JUST REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}, Di He^{1*}, Chen Dan², Kun He⁴, John E. Hopcroft³ & Liwei Wang¹
¹Peking University ²Carnegie Mellon University ³Cornell University
⁴Huazhong University of Science and Technology
 {zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
 cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon^{*} Aaditi Raghunathan^{*} Ludwig Schmidt
 Stanford University Stanford University UC Berkeley
 yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
 Stanford University Stanford University
 pliang@cs.stanford.edu jduchi@stanford.edu

Improving Robust generalisation With unlabelled data

Understanding and Mitigating the Tradeoff Between Robustness and Accuracy

Aaditi Raghunathan^{*1} Sang Michael Xie^{*1} Fanny Yang² John C. Duchi¹ Percy Liang¹

Method	Robust Test Acc.	Standard Test Acc.	
Standard Training	0.8%	95.2%	} Vanilla Supervised
PG-AT (Madry et al., 2018)	45.8%	87.3%	
TRADES (Zhang et al., 2019)	55.4%	84.0%	
Standard Self-Training	0.3%	96.4%	} Semisupervised with same unlabeled data
Robust Consistency Training (Carmon et al., 2019)	56.5%	83.2%	
RST + PG-AT (this paper)	58.5%	91.8%	
RST + TRADES (this paper)	63.1%	89.7%	

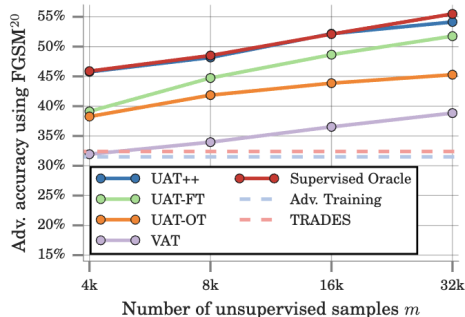
ADVERSARIALLY ROBUST GENERALIZATION JUST REQUIRES MORE UNLABELED DATA

Runtian Zhai¹; Tianle Cai^{1*}; Di He^{1*}; Chen Dan², Kun He⁴, John E. Hopcroft³ & Liwei Wang¹
¹Peking University ²Carnegie Mellon University ³Cornell University
⁴Huazhong University of Science and Technology
 {zhairuntian, caitianle1998, di.he, wanglw}@pku.edu.cn
 cdan@cs.cmu.edu, brooklet60@hust.edu.cn, jeh17@cornell.edu

Unlabeled Data Improves Adversarial Robustness

Yair Carmon^{*} Aaditi Raghunathan^{*} Ludwig Schmidt
 Stanford University Stanford University UC Berkeley
 yairc@stanford.edu aditir@stanford.edu ludwig@berkeley.edu

Percy Liang John C. Duchi
 Stanford University Stanford University
 pliang@cs.stanford.edu jduchi@stanford.edu



Are Labels Required for Improving Adversarial Robustness?

Jonathan Uesato^{*} Jean-Baptiste Alayrac^{*} Po-Sen Huang^{*}

Robert Stanforth Alhussein Fawzi Pushmeet Kohli

DeepMind
 {juesato, jalayrac, posenhuang}@google.com

Distributional Robustness in Machine Learning

Robustness to Distribution Shift

Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.

Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to

Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.

Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time

Robustness to Distribution Shift

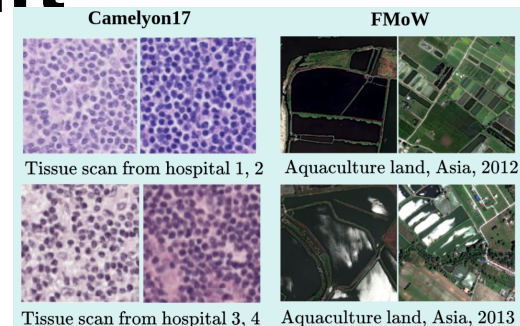
- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time

Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time

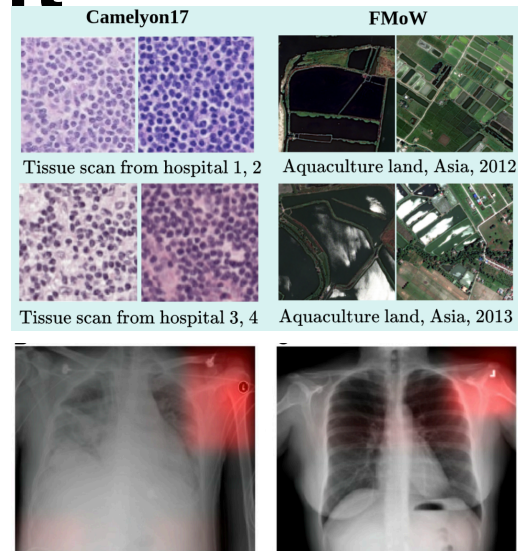
Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time



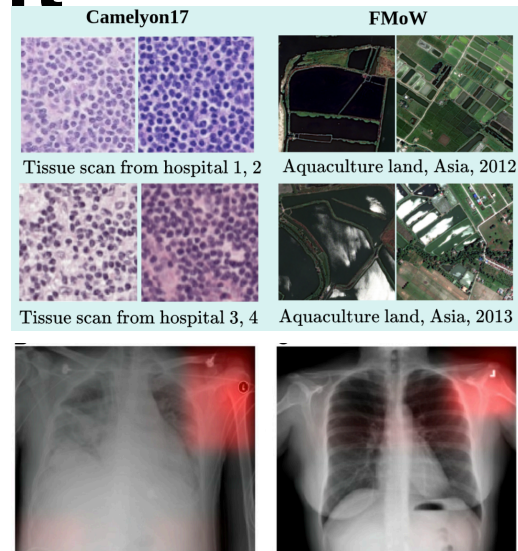
Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time



Robustness to Distribution Shift

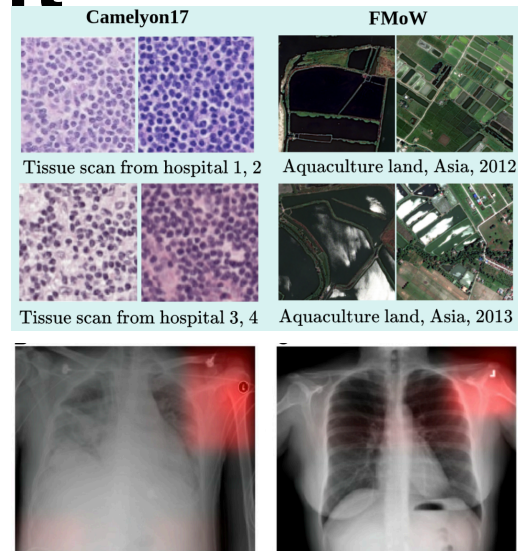
- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time



- Robustness to distribution shift requires ***preserving accuracy when the distribution shifts***

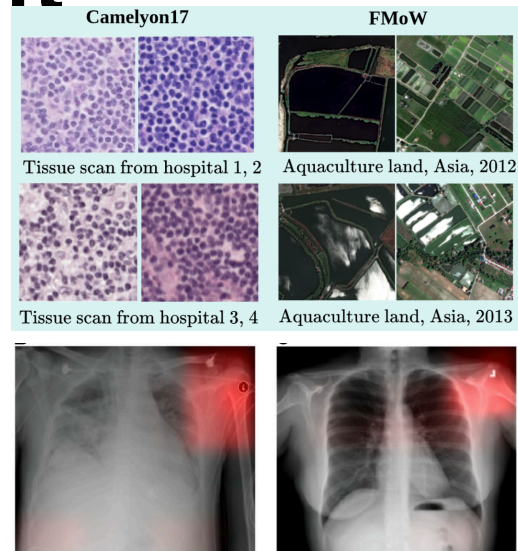
Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time
- Robustness to distribution shift requires *preserving accuracy when the distribution shifts*
- Impossible to protect against arbitrary shifts



Robustness to Distribution Shift

- **Adversarial Robustness** measures performance against the worst shift between train and test set.
- More natural distribution shifts exist in the real world between train and test data e.g. due to
 - Hard to collect quality data uniformly.
 - Data sources evolve over time



- Robustness to distribution shift requires **preserving accuracy when the distribution shifts**
- Impossible to protect against arbitrary shifts
- Goal is to allow for a graceful degradation with increasing shift

Robustness to Distribution Shift

Rich body of existing literature

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models
- **How can models be made more robust to distribution shift ?**

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models
- **How can models be made more robust to distribution shift ?**
 - Distributionally Robust Optimisation

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models
- **How can models be made more robust to distribution shift ?**
 - Distributionally Robust Optimisation
 - Learn Causal/Robust representations

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models
- **How can models be made more robust to distribution shift ?**
 - Distributionally Robust Optimisation
 - Learn Causal/Robust representations
 - Collect more data (possibly unlabelled)

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models
- **How can models be made more robust to distribution shift ?**
 - Distributionally Robust Optimisation
 - Learn Causal/Robust representations
 - Collect more data (possibly unlabelled)

Robustness to Distribution Shift

Rich body of existing literature

We will not even attempt to be exhaustive

Robustness to distribution shift features a rich body of existing literature asking

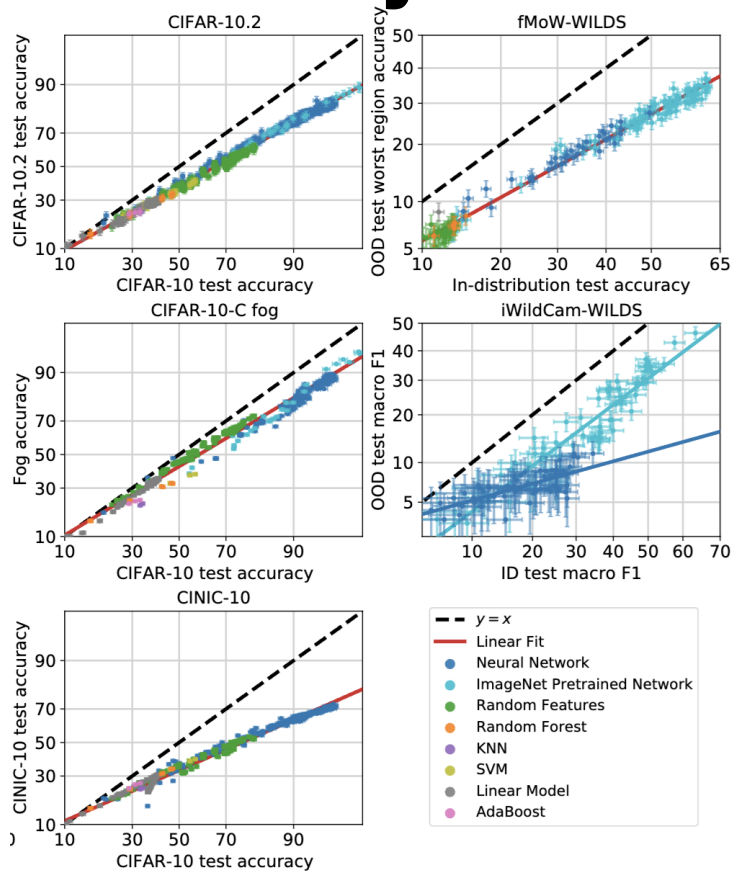
- **What causes failure to generalise to distribution shift ?**
 - Spurious correlations
 - Label Noise
 - Over-parameterised models
- **How can models be made more robust to distribution shift ?**
 - Distributionally Robust Optimisation
 - Learn Causal/Robust representations
 - Collect more data (possibly unlabelled)

Accuracy-on-the line

Accuracy on the Line: On the Strong Correlation
Between Out-of-Distribution and In-Distribution Generalization

John Miller* Rohan Taori† Aditi Raghunathan†
Shiori Sagawa† Pang Wei Koh† Vaishaal Shankar* Percy Liang†
Yair Carmon‡ Ludwig Schmidt§

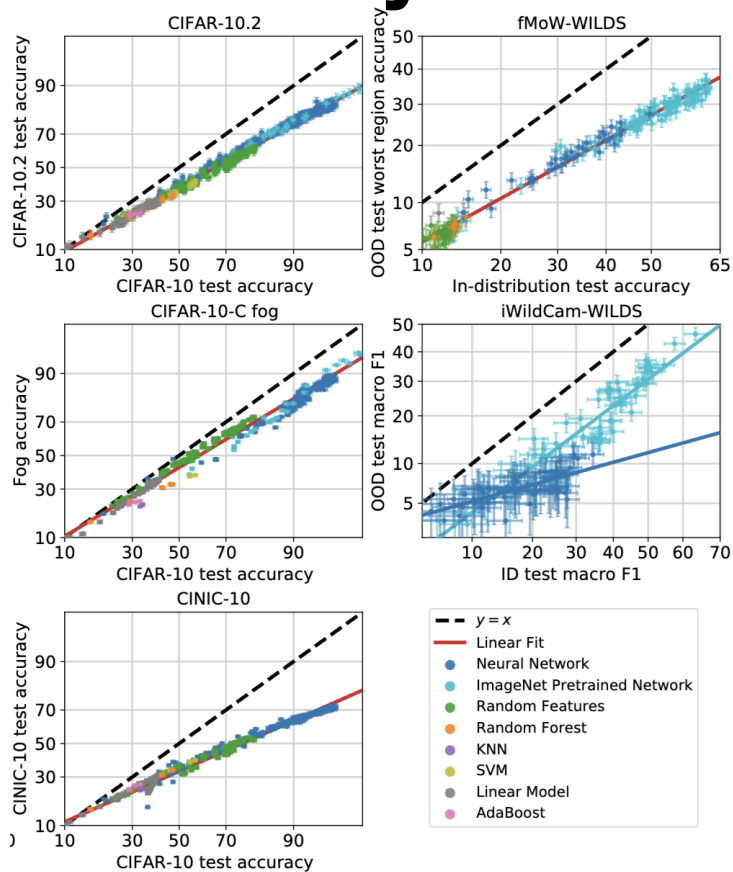
Accuracy-on-the line



Accuracy on the Line: On the Strong Correlation Between Out-of-Distribution and In-Distribution Generalization

John Miller* Rohan Taori† Aditi Raghunathan†
 Shiori Sagawa† Pang Wei Koh† Vaishaal Shankar* Percy Liang†
 Yair Carmon‡ Ludwig Schmidt§

Accuracy-on-the line

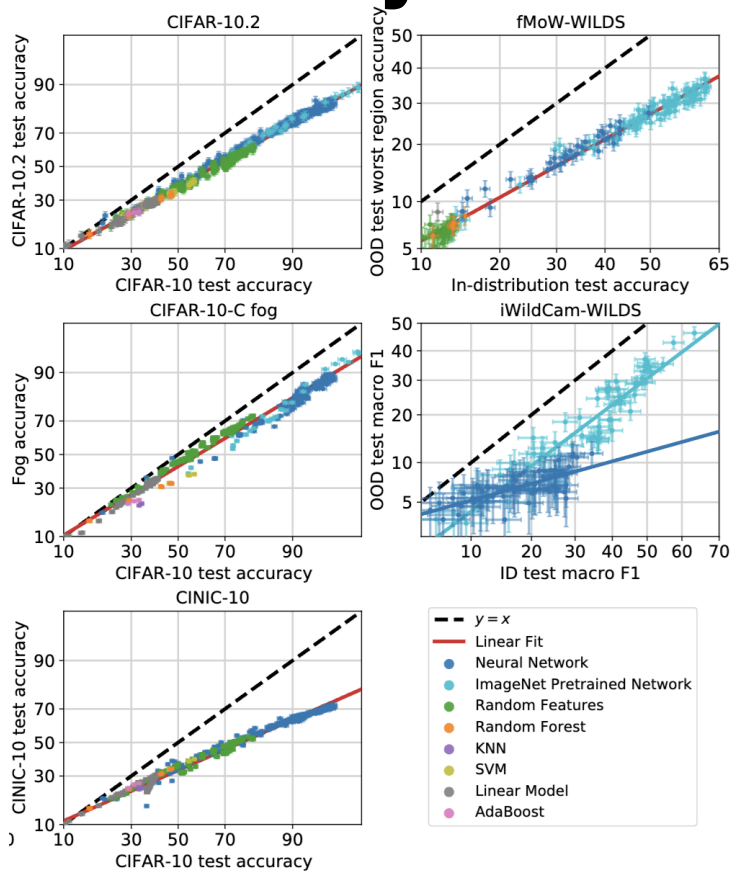


Accuracy on the Line: On the Strong Correlation
Between Out-of-Distribution and In-Distribution Generalization

John Miller* Rohan Taori† Aditi Raghunathan†
Shiori Sagawa† Pang Wei Koh† Vaishaal Shankar* Percy Liang†
Yair Carmon‡ Ludwig Schmidt§

- **Accuracy-on-the-line** phenomenon: ID and OOD accuracy are positively correlated.

Accuracy-on-the line

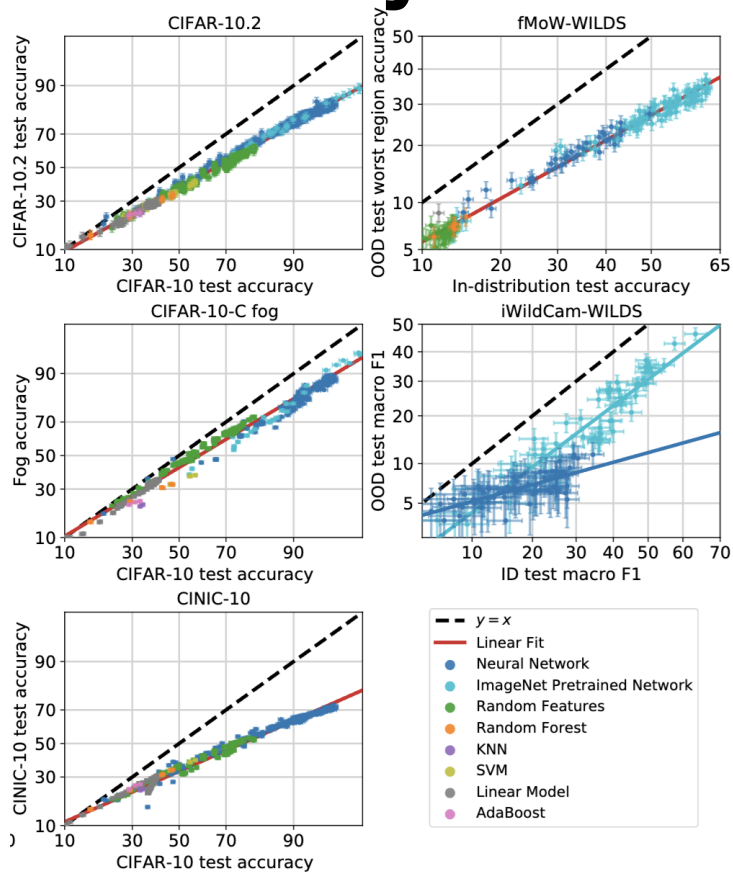


Accuracy on the Line: On the Strong Correlation
Between Out-of-Distribution and In-Distribution Generalization

John Miller* Rohan Taori† Aditi Raghunathan†
Shiori Sagawa† Pang Wei Koh† Vaishal Shankar* Percy Liang†
Yair Carmon‡ Ludwig Schmidt§

- **Accuracy-on-the-line** phenomenon: ID and OOD accuracy are positively correlated.
- Indicates that improving ID accuracy also improves OOD accuracy.

Accuracy-on-the line



Accuracy on the Line: On the Strong Correlation
Between Out-of-Distribution and In-Distribution Generalization

John Miller* Rohan Taori† Aditi Raghunathan†
Shiori Sagawa† Pang Wei Koh† Vaishaal Shankar* Percy Liang†
Yair Carmon‡ Ludwig Schmidt§

- **Accuracy-on-the-line** phenomenon: ID and OOD accuracy are positively correlated.
- Indicates that improving ID accuracy also improves OOD accuracy.
- Holds for a wide variety of models and datasets

Label Noise and Distribution Shift

Accuracy on the wrong line: On the pitfalls of noisy data for
out-of-distribution generalisation

Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

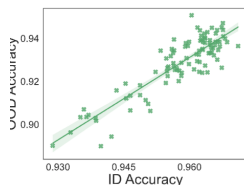
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.



(b) Noiseless dataset

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

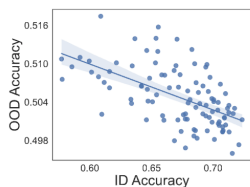
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

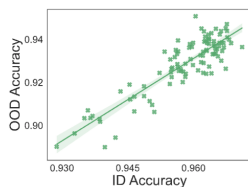
²University of California, Berkeley, U.S.A.

³Halcioğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.



(a) Noisy dataset



(b) Noiseless dataset

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

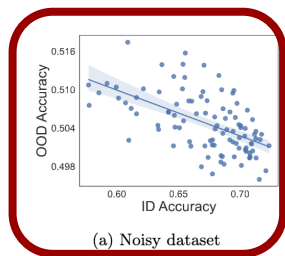
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

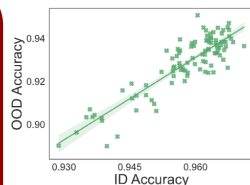
²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.



(a) Noisy dataset



(b) Noiseless dataset

Accuracy-on-the-wrong line

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

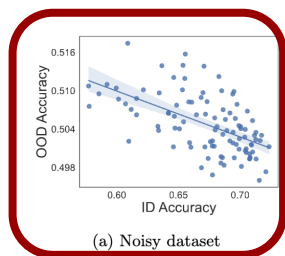
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

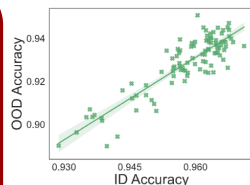
²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

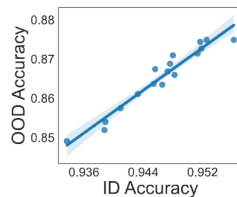
⁴University of Michigan, Ann Arbor, U.S.A.



(a) Noisy dataset



(b) Noiseless dataset



(a) $\eta = 0$.

Accuracy-on-the-wrong line

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

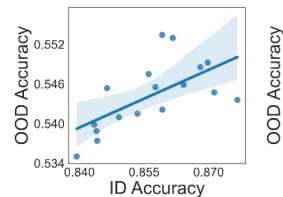
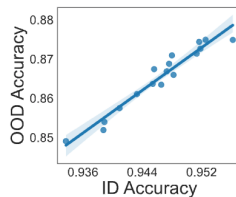
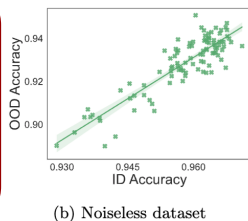
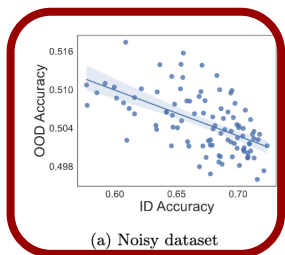
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.



Accuracy-on-the-wrong line

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

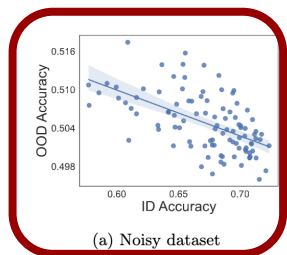
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

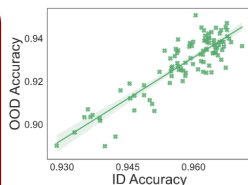
²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

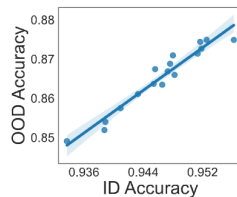
⁴University of Michigan, Ann Arbor, U.S.A.



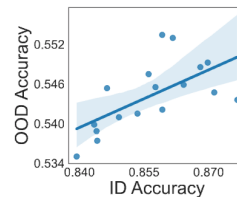
(a) Noisy dataset



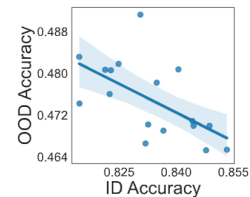
(b) Noiseless dataset



(a) $\eta = 0$.



(b) $\eta = 0.15$



(c) $\eta = 0.2$

Accuracy-on-the-wrong line

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

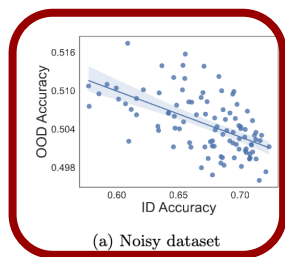
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

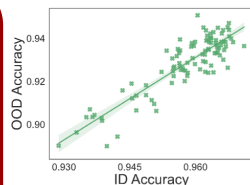
²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.

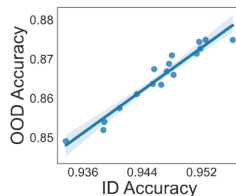


(a) Noisy dataset

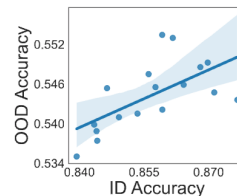


(b) Noiseless dataset

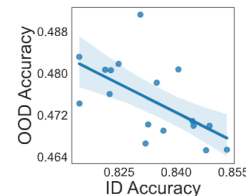
Accuracy-on-the-wrong line



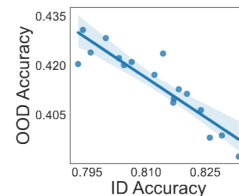
(a) $\eta = 0$.



(b) $\eta = 0.15$



(c) $\eta = 0.2$



(d) $\eta = 0.25$

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

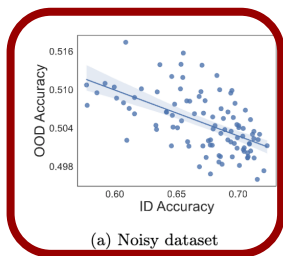
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

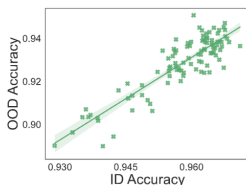
²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

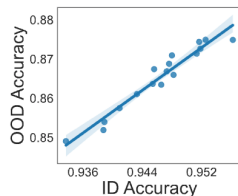
⁴University of Michigan, Ann Arbor, U.S.A.



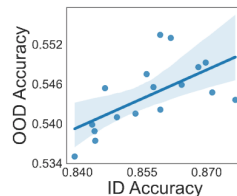
(a) Noisy dataset



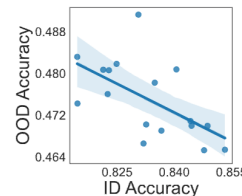
(b) Noiseless dataset



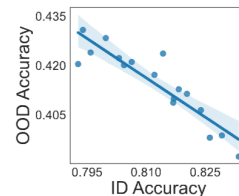
(a) $\eta = 0$.



(b) $\eta = 0.15$



(c) $\eta = 0.2$



(d) $\eta = 0.25$

Accuracy-on-the-wrong line

Two sufficient factors for Accuracy-on-the-wrong-line

Label Noise and Distribution Shift

- Question: Is **Accuracy-on-the-line** robust to noisy or low quality labels ?

Accuracy on the wrong line: On the pitfalls of noisy data for out-of-distribution generalisation

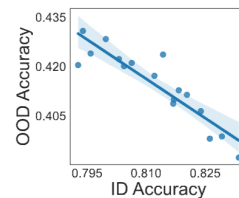
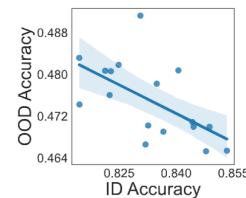
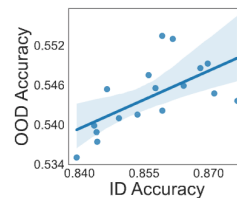
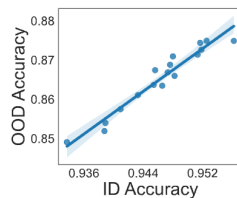
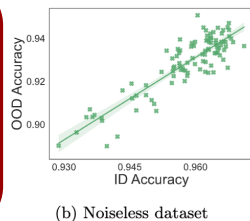
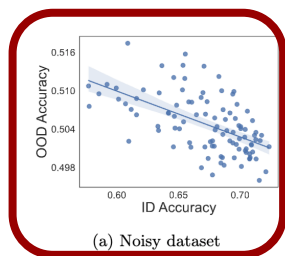
Amartya Sanyal¹, Yaxi Hu¹, Yaodong Yu², Yian Ma³, Yixin Wang⁴, and Bernhard Schölkopf¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

²University of California, Berkeley, U.S.A.

³Hacıoğlu Data Science Institute, University of California San Diego, San Diego, U.S.A.

⁴University of Michigan, Ann Arbor, U.S.A.



Accuracy-on-the-wrong line

Two sufficient factors for Accuracy-on-the-wrong-line

- Inject and fit **random label noise** in the training data
- Presence of multiple **“nuisance features”** i.e. irrelevant features

Theory

Theory

Let the data satisfy the following

Theory

Let the data satisfy the following



Theory

Let the data satisfy the following



Signal support S_0

Theory

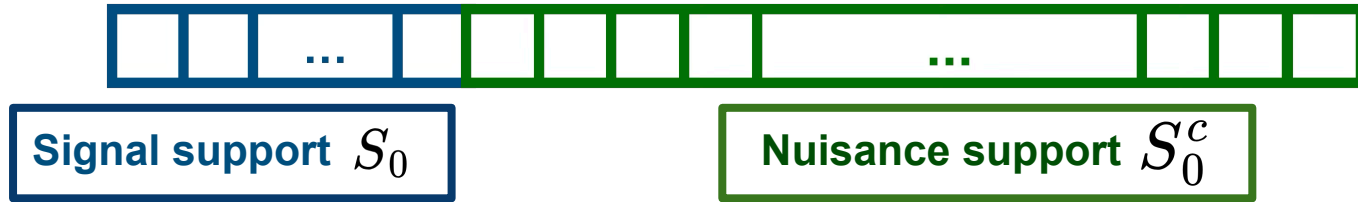
Let the data satisfy the following



Signal support S_0

Theory

Let the data satisfy the following



Theory

Let the data satisfy the following



Signal support S_0

Nuisance support S_0^c

Data is labelled as $y = \langle \theta^*, x \rangle$
where θ^* is supported on S_0

Theory

Let the data satisfy the following



Signal support S_0

Nuisance support S_0^c

Data is labelled as $y = \langle \theta^*, x \rangle$
where θ^* is supported on S_0

S_0 -oblivious shift distribution Δ only
affects nuisance support S_0^c .

Theory

Let the data satisfy the following



Signal support S_0

Nuisance support S_0^c

Data is labelled as $y = \langle \theta^*, x \rangle$
where θ^* is supported on S_0

S_0 -oblivious shift distribution Δ only
affects nuisance support S_0^c .

Implication of label noise: Learned linear classifier has spurious sensitivity $\geq \tau$

Theory

Let the data satisfy the following



Signal support S_0

Nuisance support S_0^c

Data is labelled as $y = \langle \theta^*, x \rangle$
where θ^* is supported on S_0

S_0 -oblivious shift distribution Δ only
affects nuisance support S_0^c .

Implication of label noise: Learned linear classifier has spurious sensitivity $\geq \tau$

Informal Theorem For all x s.t. $\langle \theta^*, x \rangle > 0$, we have

Theory

Let the data satisfy the following



Signal support S_0

Nuisance support S_0^c

Data is labelled as $y = \langle \theta^*, x \rangle$
where θ^* is supported on S_0

S_0 -oblivious shift distribution Δ only
affects nuisance support S_0^c .

Implication of label noise: Learned linear classifier has spurious sensitivity $\geq \tau$

Informal Theorem For all x s.t. $\langle \theta^*, x \rangle > 0$, we have

$$\Pr_{\delta \sim \Delta} \left[\langle \hat{\theta}, x + \delta \rangle < 0 \right] \geq 1 - \exp \left(- |S_0^c| \tau^2 \right)$$

Using unlabelled data

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*

Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt

Department of Computer Science
ETH Zurich

Philip H.S. Torr

Department of Engineering Science
University of Oxford

Amartya Sanyal

Department of Computer Science & ETH AI Center
ETH Zurich

Using unlabelled data

- **Pre-trained representations** are a common strategy against this problem.
- But representations from supervised training often suffer from problems like **simplicity bias**.

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

Using unlabelled data

- **Pre-trained representations** are a common strategy against this problem.
- But representations from supervised training often suffer from problems like **simplicity bias**.

Solution - Use **Unlabelled data** & **unsupervised representation** learning

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

Using unlabelled data

- **Pre-trained representations** are a common strategy against this problem.
- But representations from supervised training often suffer from problems like **simplicity bias**.

Solution - Use **Unlabelled data** & **unsupervised representation** learning

Experimental setup

HOW ROBUST IS UNSUPERVISED REPRESENTATION
LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

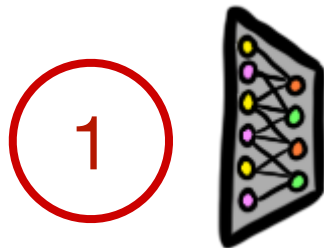
Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

Using unlabelled data

- **Pre-trained representations** are a common strategy against this problem.
- But representations from supervised training often suffer from problems like **simplicity bias**.

Solution - Use **Unlabelled data** & **unsupervised representation** learning

Experimental setup



Pre-train representation learning on **ID data**
with labelled (SL) or unlabelled data (AE/SSL)

HOW ROBUST IS UNSUPERVISED REPRESENTATION
LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

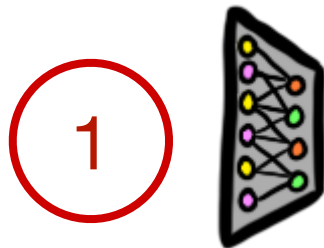
Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

Using unlabelled data

- **Pre-trained representations** are a common strategy against this problem.
- But representations from supervised training often suffer from problems like **simplicity bias**.

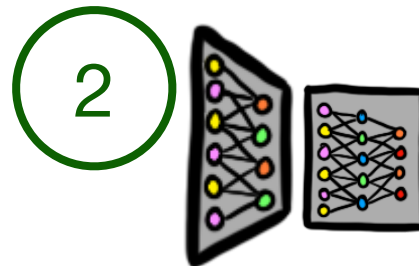
Solution - Use **Unlabelled data** & **unsupervised representation** learning

Experimental setup



Pre-train representation learning on **ID data** with labelled (SL) or unlabelled data (AE/SSL)

Train a small ML model on top of the features using **Dist X (ID or OOD)**



HOW ROBUST IS UNSUPERVISED REPRESENTATION
LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

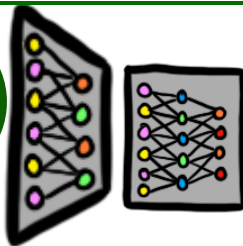
Using unlabelled data

1



Train a small ML model on top of the features using **Dist X**

2



Pre-train representation learning on **ID data** with labelled (SL) or unlabelled data (AE/SSL)

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

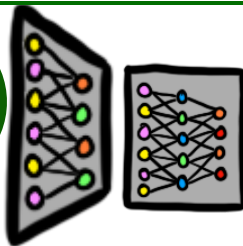
Using unlabelled data

1



Train a small ML model on top of the features using **Dist X**

2



Pre-train representation learning on **ID data** with labelled (SL) or unlabelled data (AE/SSL)

Dist X \rightarrow OOD. Test on OOD.

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

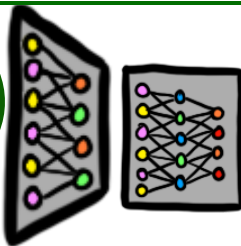
Using unlabelled data

1



Train a small ML model on top of the features using **Dist X**

2



Pre-train representation learning on **ID data** with labelled (SL) or unlabelled data (AE/SSL)

Dist X \rightarrow OOD. Test on OOD.

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

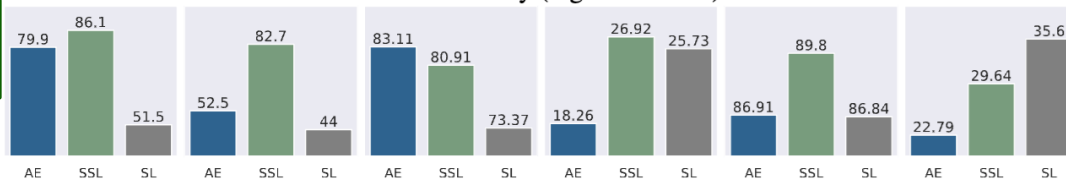
Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

OOD Accuracy (higher is better)



(a) MNIST-CIFAR

(b) CdSprites

(c) Camelyon17-CS

(d) FMoW-CS

(e) Camelyon17

(f) FMoW

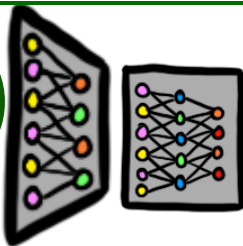
Using unlabelled data

1



Train a small ML model on top of the features using **Dist X**

2



Pre-train representation learning on **ID data** with labelled (SL) or unlabelled data (AE/SSL)

Dist X \rightarrow OOD. Test on OOD.

Shift Sensitivity = Diff between

1. Dist X \rightarrow OOD. Test on OOD.
2. Dist X \rightarrow ID. Test on ID.

Captures robustness of

HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

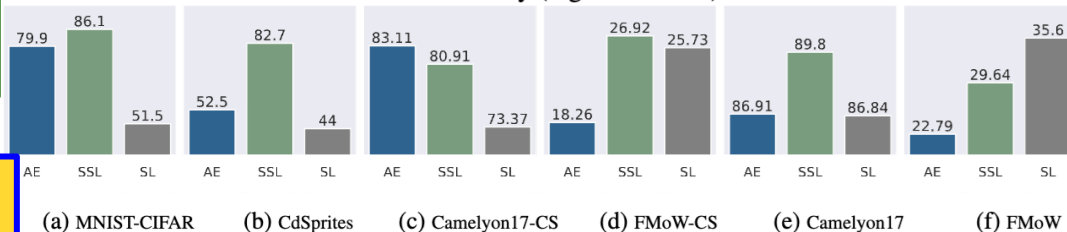
Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

OOD Accuracy (higher is better)



(a) MNIST-CIFAR

(b) CdSprites

(c) Camelyon17-CS

(d) FMoW-CS

(e) Camelyon17

(f) FMoW

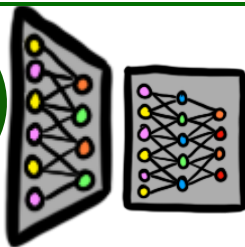
Using unlabelled data

1



Train a small ML model on top of the features using **Dist X**

2



Pre-train representation learning on **ID data** with labelled (SL) or unlabelled data (AE/SSL)

Dist X \rightarrow OOD. Test on OOD.

Shift Sensitivity = Diff between

- Dist X \rightarrow OOD. Test on OOD.
- Dist X \rightarrow ID. Test on ID.

Captures robustness of

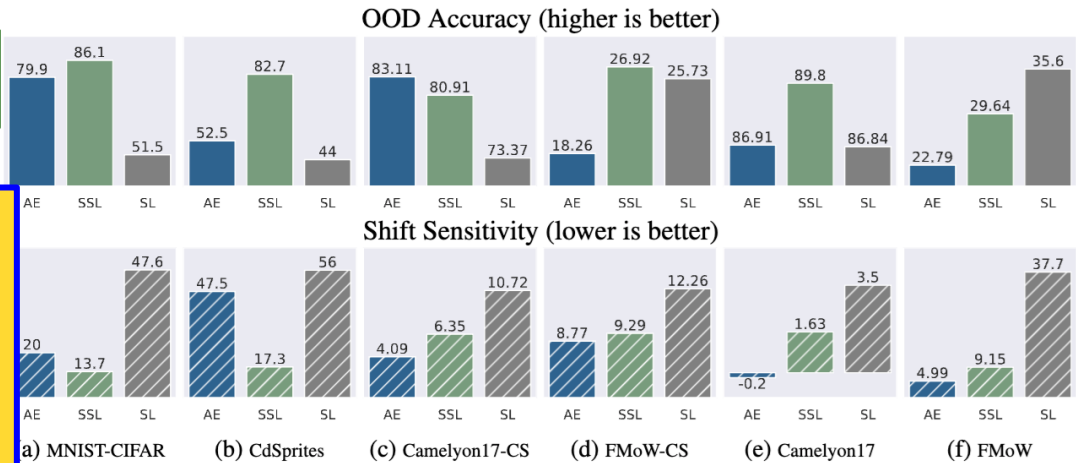
HOW ROBUST IS UNSUPERVISED REPRESENTATION LEARNING TO DISTRIBUTION SHIFT?

Yuge Shi*
Department of Engineering Science
University of Oxford

Imant Daunhawer & Julia E. Vogt
Department of Computer Science
ETH Zurich

Philip H.S. Torr
Department of Engineering Science
University of Oxford

Amartya Sanyal
Department of Computer Science & ETH AI Center
ETH Zurich

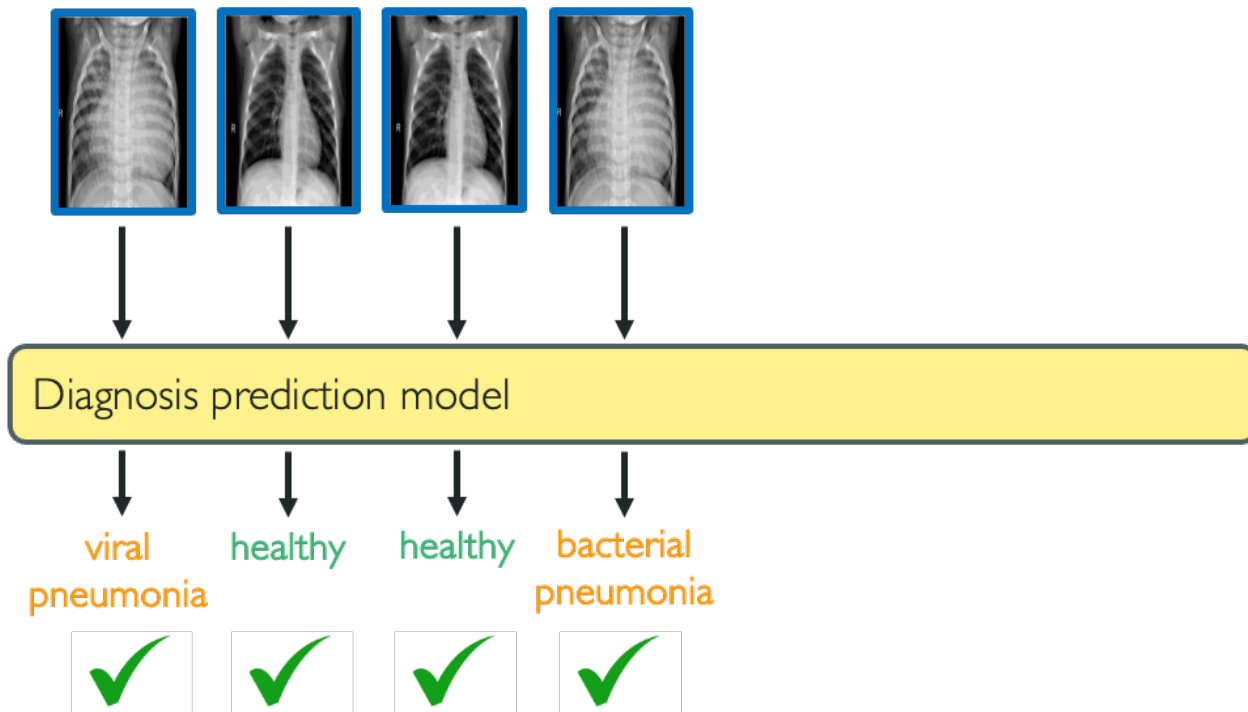


Out-of-distribution detection

**What if we cannot predict reliably
outside of the training distribution?**

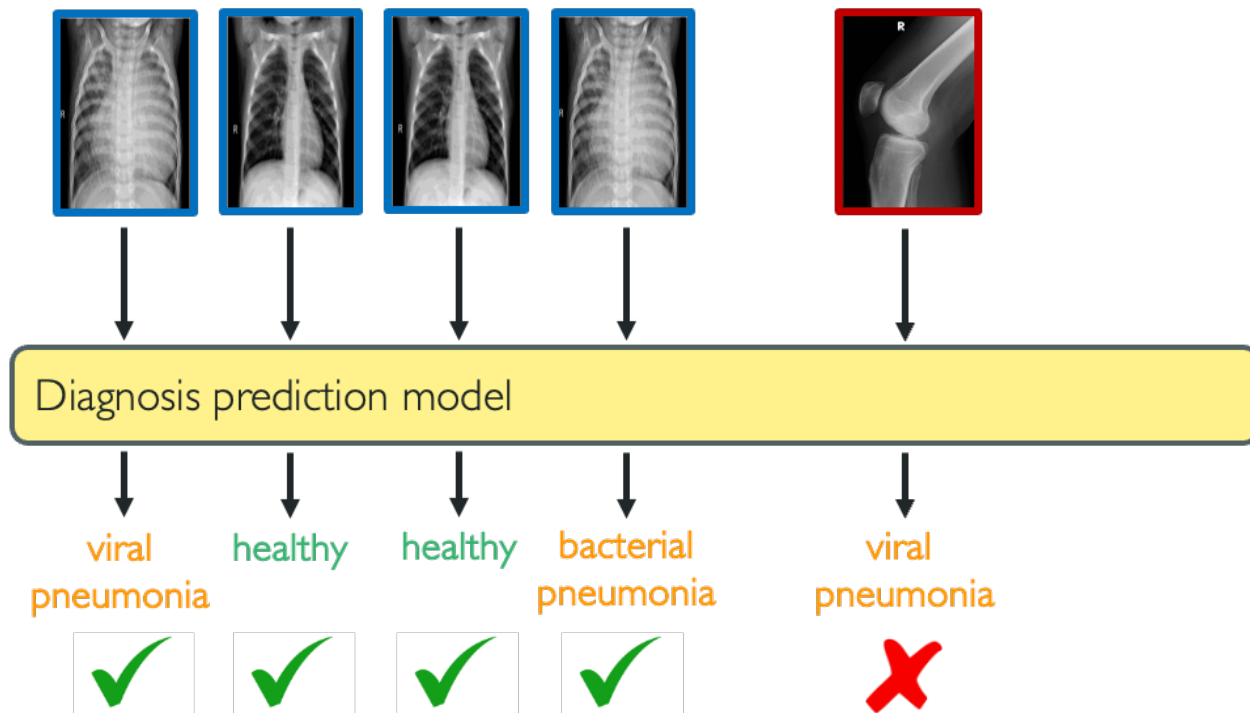
When can't we predict on OOD data?

Novel classes



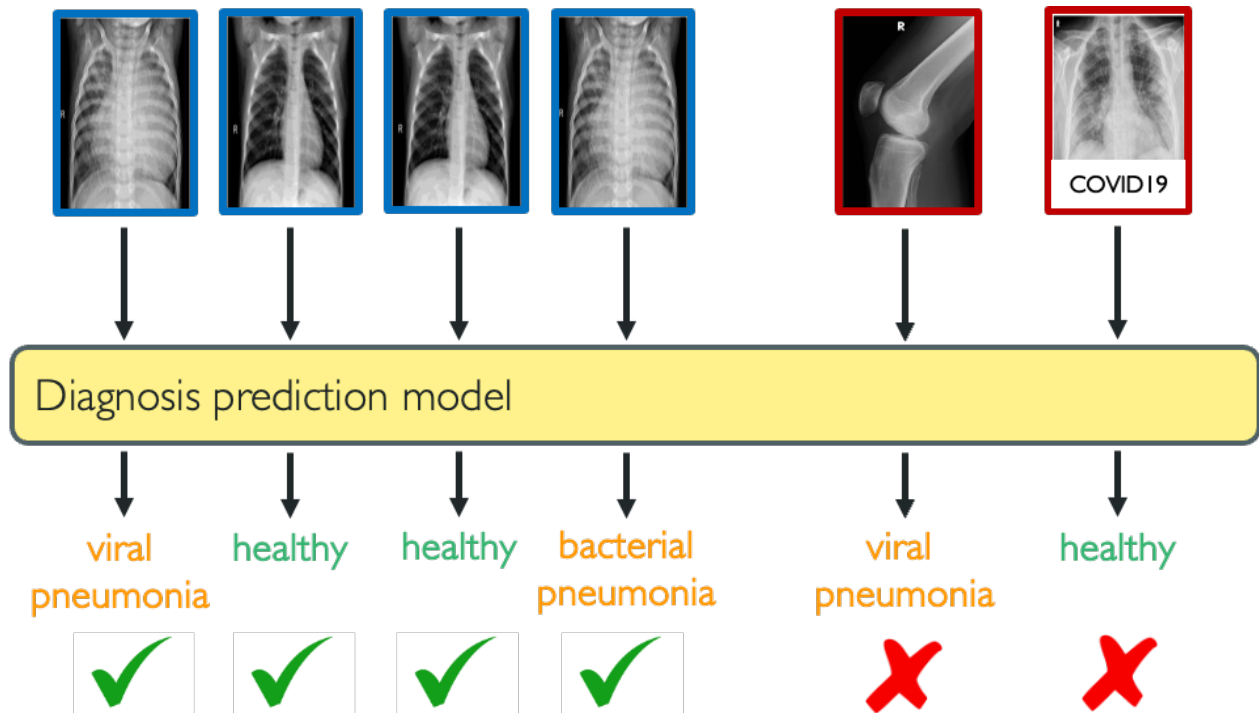
When can't we predict on OOD data?

Novel classes



When can't we predict on OOD data?

Novel classes



When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

invariant
parameters

domain-specific
parameters

When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

invariant
parameters

domain-specific
parameters

θ^*



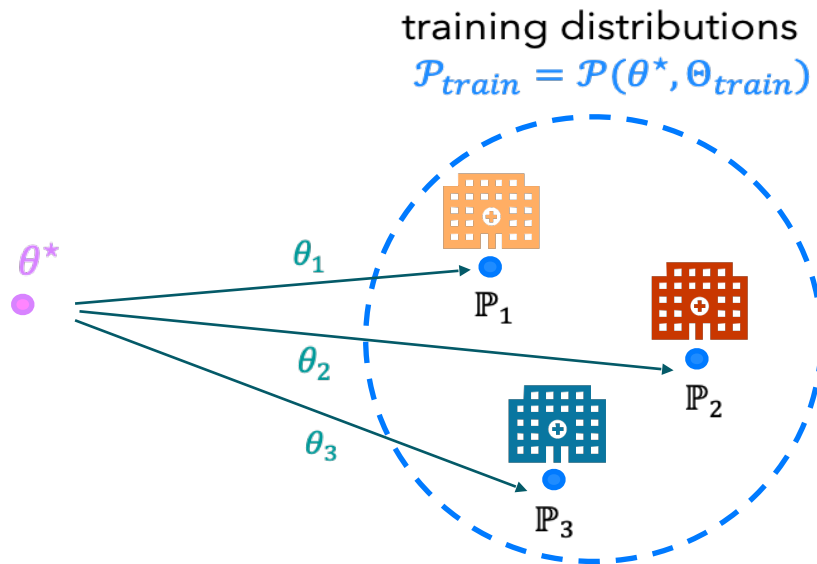
When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

invariant
parameters

domain-specific
parameters



When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

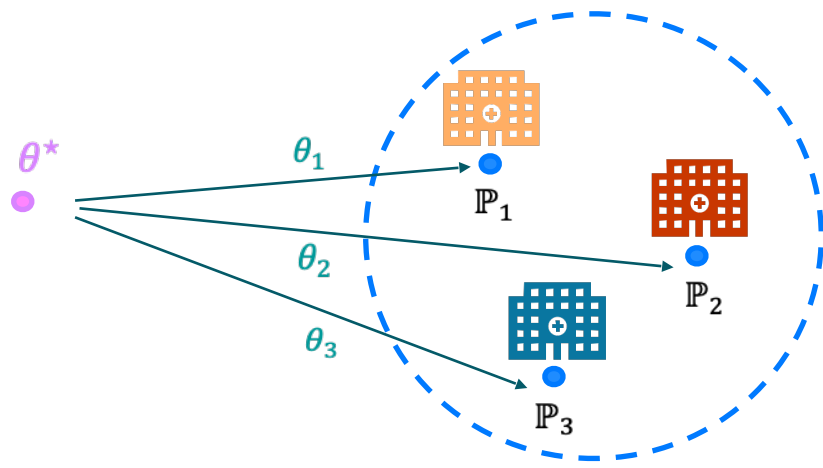
invariant
parameters

domain-specific
parameters

Model
 β



training distributions
 $\mathcal{P}_{train} = \mathcal{P}(\theta^*, \Theta_{train})$



When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

invariant
parameters

domain-specific
parameters

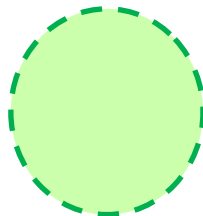
Model
 β



training distributions
 $\mathcal{P}_{train} = \mathcal{P}(\theta^*, \Theta_{train})$

possible test distributions

$$\mathcal{P}_{test} = \mathcal{P}(\theta^*, \Theta_{test})$$



possible
set Θ_{test} of
test shifts

θ^*

θ_1

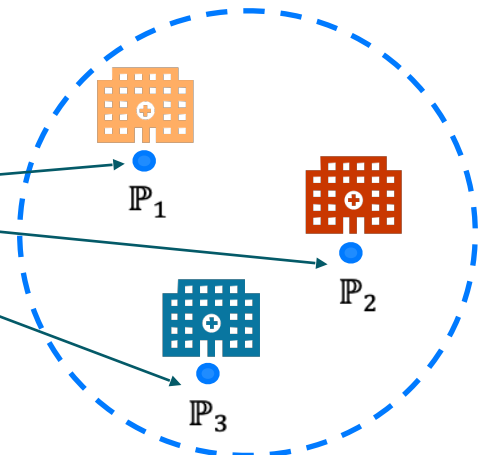
θ_2

θ_3

\mathbb{P}_1

\mathbb{P}_2

\mathbb{P}_3



When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

invariant parameters

domain-specific parameters

Model β



training distributions
 $\mathcal{P}_{train} = \mathcal{P}(\theta^*, \Theta_{train})$

possible test distributions

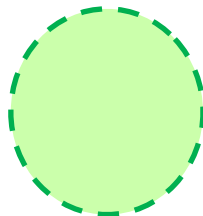
$$\mathcal{P}_{test} = \mathcal{P}(\theta^*, \Theta_{test})$$

robust risk

$$R_{rob}(\beta; \theta^*, \Theta_{test})$$

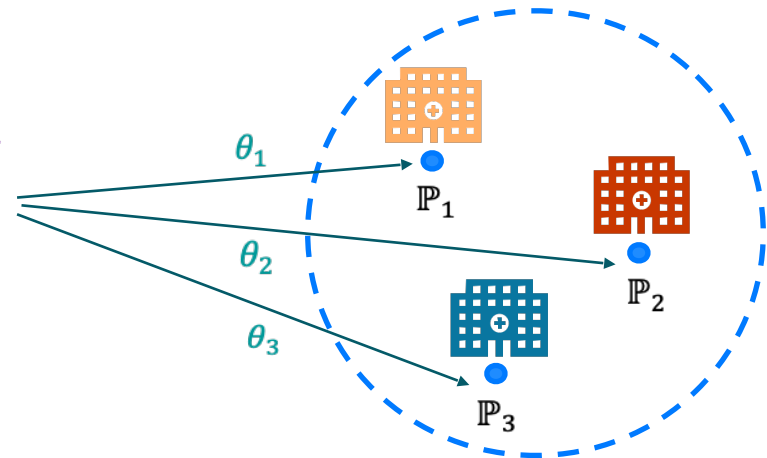


worst-case risk in \mathcal{P}_{test} for β



possible set Θ_{test} of test shifts

θ^*



When can't we predict on OOD data?

Strong distribution shifts

$\mathbb{P}(X, Y)$ determined by (θ^*, θ_e)

invariant parameters

domain-specific parameters

Model β



training distributions
 $\mathcal{P}_{train} = \mathcal{P}(\theta^*, \Theta_{train})$

possible test distributions

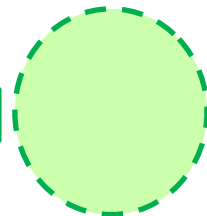
$$\mathcal{P}_{test} = \mathcal{P}(\theta^*, \Theta_{test})$$

robust risk

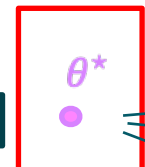
$$R_{rob}(\beta; \theta^*, \Theta_{test})$$



worst-case risk in \mathcal{P}_{test} for β



possible set Θ_{test} of test shifts

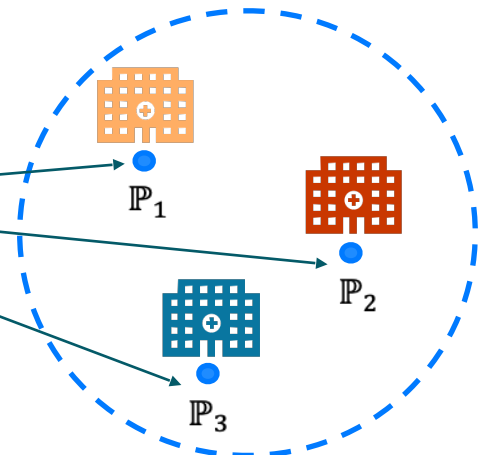


not generally identifiable!

θ_1

θ_2

θ_3



Impossibility result for distribution shifts

Achievable distributional robustness when the robust risk is only partially identified

Julia Kostin¹, Nicola Gnecco², and Fanny Yang¹

¹Department of Computer Science, ETH Zurich

²Department of Mathematics, Imperial College London

Mean shifts during test time assumed to lie in $\Theta_{test} = \{\theta_{test}: \theta_{test}\theta_{test}^\top \preceq \gamma M_{seen} + \gamma' M_{unseen}\}$

Test time shifts assumptions

Covariance with range in span of **seen shift directions** $range(M_{seen}) \subset span\{\theta_e\}_{e \in [k]}$

Projection matrix onto **unseen direction**: $range(M_{seen}) \perp span\{\theta_e\}_{e \in [k]}$

Impossibility result for distribution shifts

Achievable distributional robustness when the robust risk is only partially identified

Julia Kostin¹, Nicola Gnecco², and Fanny Yang¹

¹Department of Computer Science, ETH Zurich

²Department of Mathematics, Imperial College London

Mean shifts during test time assumed to lie in $\Theta_{test} = \{\theta_{test}: \theta_{test}\theta_{test}^\top \preceq \gamma M_{seen} + \gamma' M_{unseen}\}$
shift strengths

Test time shifts assumptions

Covariance with range in span of **seen shift directions** $range(M_{seen}) \subset span\{\theta_e\}_{e \in [k]}$

Projection matrix onto **unseen direction**: $range(M_{seen}) \perp span\{\theta_e\}_{e \in [k]}$

Impossibility result for distribution shifts

Achievable distributional robustness when the robust risk is only partially identified

Julia Kostin¹, Nicola Gnecco², and Fanny Yang¹

¹Department of Computer Science, ETH Zurich
²Department of Mathematics, Imperial College London

Mean shifts during test time assumed to lie in $\Theta_{test} = \{\theta_{test}: \theta_{test}\theta_{test}^\top \preceq \gamma M_{seen} + \gamma' M_{unseen}\}$
shift strengths

Test time shifts assumptions

Covariance with range in span of **seen shift directions** $range(M_{seen}) \subset span\{\theta_e\}_{e \in [k]}$

Projection matrix onto **unseen direction**: $range(M_{seen}) \perp span\{\theta_e\}_{e \in [k]}$

Main theoretical result

Information-theoretic lower bound on robust risk.

Corollary

- **No “unseen” shifts:** Existing OOD generalization algorithms (e.g. anchor regression) are optimal.
- **No “seen” shifts:** Anchor regression is not better than ordinary least squares.

**What if we cannot predict reliably
outside of the training distribution?**

**What if we cannot predict reliably
outside of the training distribution?**

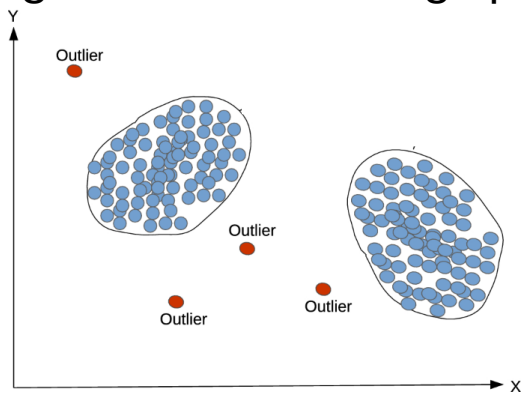
**A: Flag out-of-domain samples
and abstain.**

Traditional OOD detection methods

Unsupervised OOD i.e. only observe in-distribution samples.

Examples:

Density estimation
e.g. in NN embedding space



Predictive uncertainty
e.g. ensembles

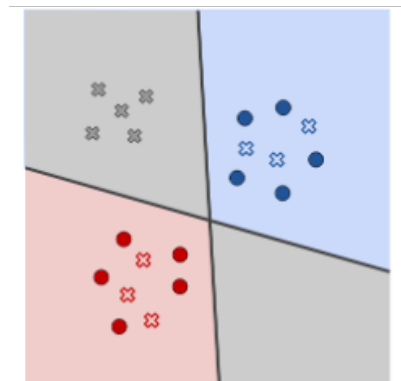


Figure sources: <https://link.springer.com/article/10.1007/s10044-021-00998-6>, <https://arxiv.org/abs/2012.05825>

Limitations of unsupervised OOD detection

Limitations of unsupervised OOD detection

Challenge #1: Unsupervised OOD detection can be ill-defined

Perfect Density Models Cannot Guarantee Anomaly Detection

Charline Le Lan ^{1,2,*} and Laurent Dinh ²

¹ Department of Statistics, University of Oxford, Oxford OX1 3LB, UK

² Google Research, Montreal H3B 2Y5, CA

Limitations of unsupervised OOD detection

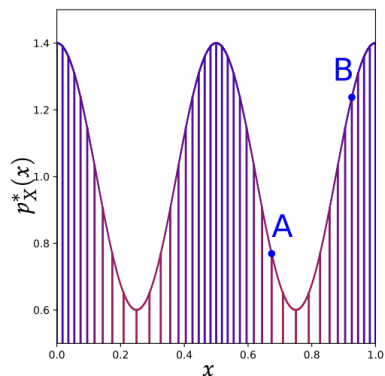
Challenge #1: Unsupervised OOD detection can be ill-defined

Perfect Density Models Cannot Guarantee Anomaly Detection

Charline Le Lan ^{1,2,*} and Laurent Dinh ²

¹ Department of Statistics, University of Oxford, Oxford OX1 3LB, UK

² Google Research, Montreal H3B 2Y5, CA



Limitations of unsupervised OOD detection

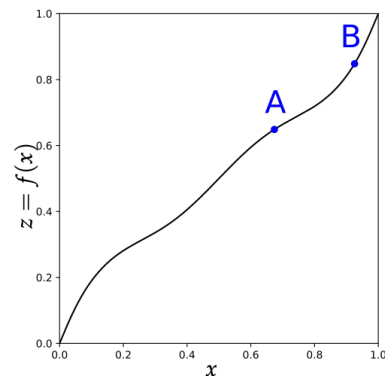
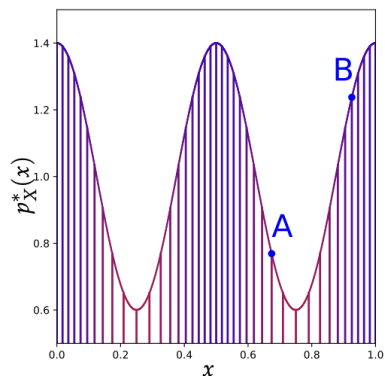
Challenge #1: Unsupervised OOD detection can be ill-defined

Perfect Density Models Cannot Guarantee Anomaly Detection

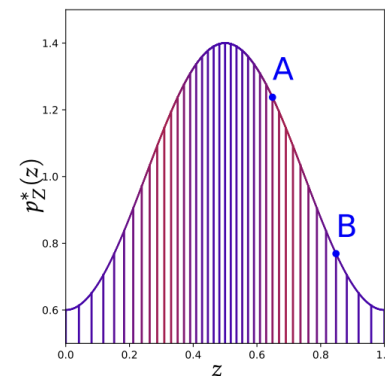
Charline Le Lan^{1,2*} and Laurent Dinh²

¹ Department of Statistics, University of Oxford, Oxford OX1 3LB, UK

² Google Research, Montreal H3B 2Y5, CA



invertible change of representation



Limitations of unsupervised OOD detection

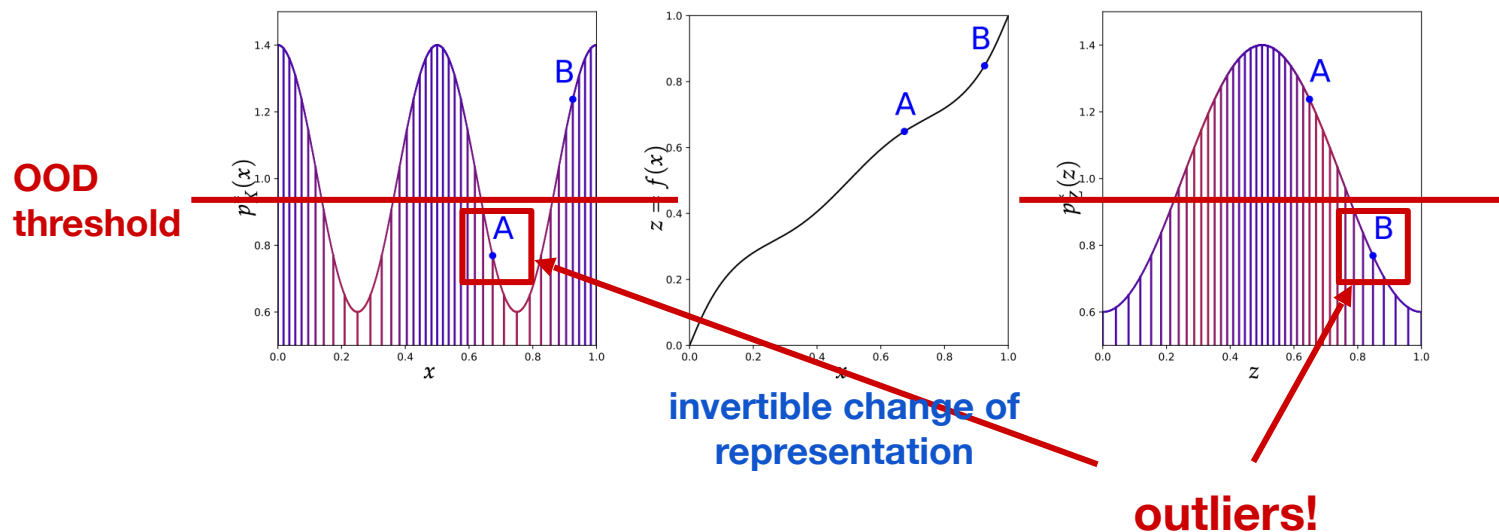
Challenge #1: Unsupervised OOD detection can be ill-defined

Perfect Density Models Cannot Guarantee Anomaly Detection

Charline Le Lan^{1,2*} and Laurent Dinh²

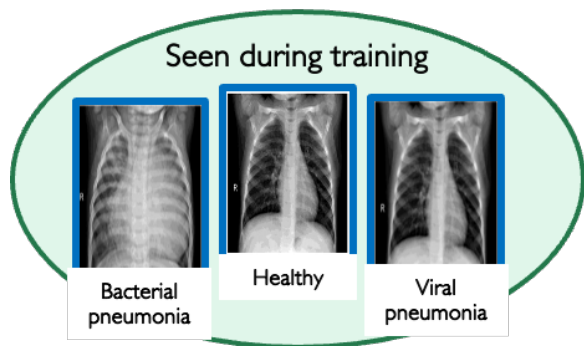
¹ Department of Statistics, University of Oxford, Oxford OX1 3LB, UK

² Google Research, Montreal H3B 2Y5, CA



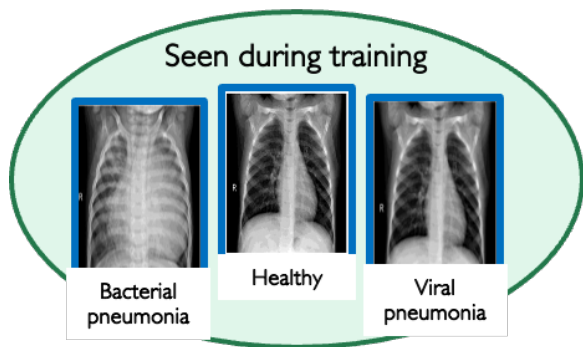
Limitations of unsupervised OOD detection

Challenge #2: finite samples + curse of dimensionality



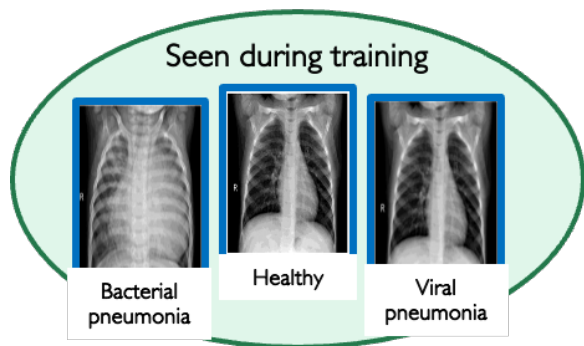
Limitations of unsupervised OOD detection

Challenge #2: finite samples + curse of dimensionality

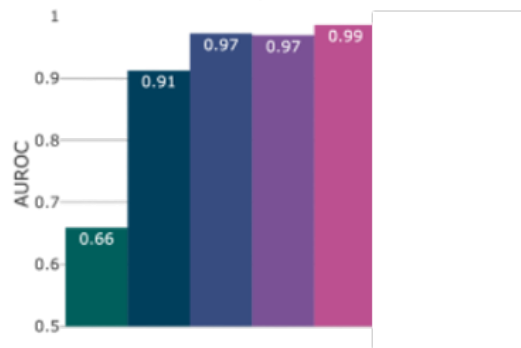


Limitations of unsupervised OOD detection

Challenge #2: finite samples + curse of dimensionality

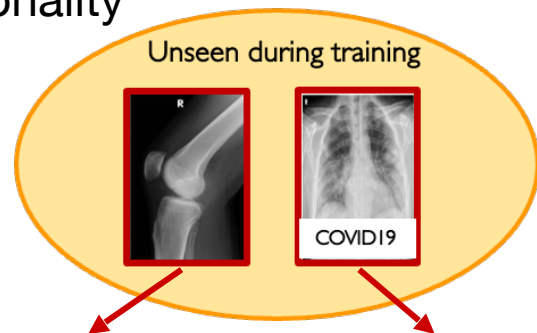
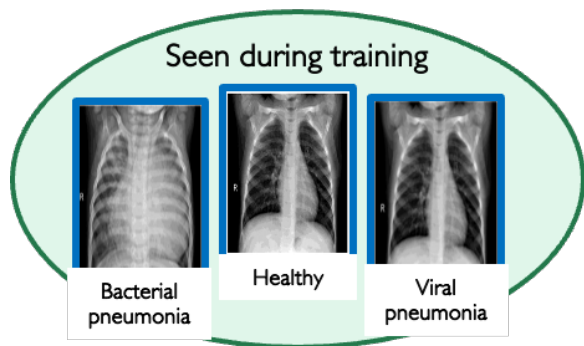


Far OOD ("easy")
i.e. OOD = new dataset

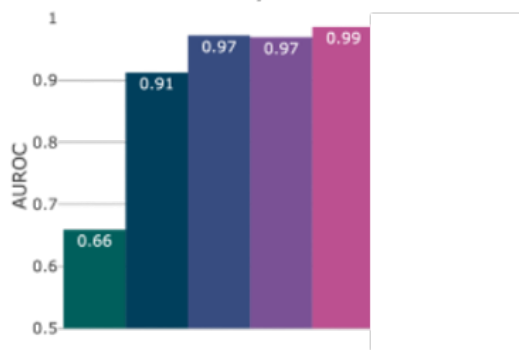


Limitations of unsupervised OOD detection

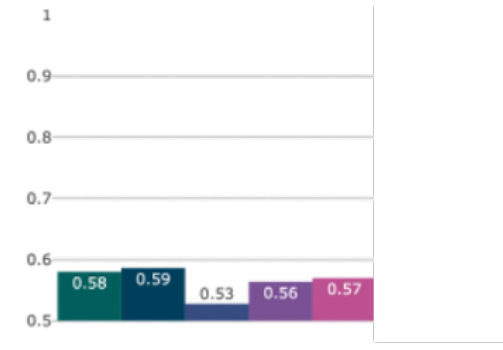
Challenge #2: finite samples + curse of dimensionality



Far OOD (“easy”)
i.e. OOD = new dataset



Near OOD (“hard”)
i.e. OOD = new classes



Diverse pre-training data

Pre-train on ImageNet21k



Exploring the Limits of Out-of-Distribution Detection

Stanislav Fort*
Stanford University
sfort1@stanford.edu

Jie Ren*
Google Research, Brain Team
jjren@google.com

Balaji Lakshminarayanan
Google Research, Brain Team
balajiln@google.com

Diverse pre-training data

Exploring the Limits of Out-of-Distribution Detection

Stanislav Fort*
Stanford University
sfort1@stanford.edu

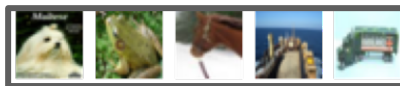
Jie Ren*
Google Research, Brain Team
jjren@google.com

Balaji Lakshminarayanan
Google Research, Brain Team
balajiln@google.com

Pre-train on ImageNet21k



Fine-tune on CIFAR10



Outliers: CIFAR100



Unsup. method:

Pretrained method:

AUROC

0.80

0.97

Diverse pre-training data

Exploring the Limits of Out-of-Distribution Detection

Stanislav Fort*
Stanford University
sfort1@stanford.edu

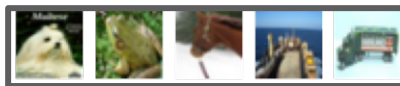
Jie Ren*
Google Research, Brain Team
jjren@google.com

Balaji Lakshminarayanan
Google Research, Brain Team
balajiln@google.com

Pre-train on ImageNet21k



Fine-tune on CIFAR10



Outliers: CIFAR100



AUROC

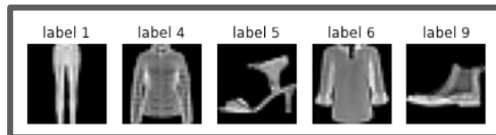
Unsup. method:

0.80

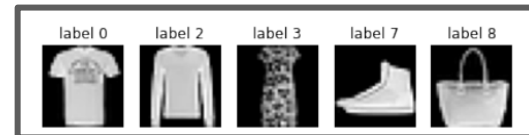
Pretrained method:

0.97

Fine-tune on
5-class FashionMNIST



Outliers: remaining
FashionMNIST classes



AUROC

Unsup. method:

0.82

Pretrained method:

0.87

Using proxy OOD data

Natural proxy OOD data

DEEP ANOMALY DETECTION WITH OUTLIER EXPOSURE

Dan Hendrycks

University of California, Berkeley
hendrycks@berkeley.edu

Mantas Mazeika

University of Chicago
mantas@ttic.edu

Thomas Dietterich

Oregon State University
tgd@oregonstate.edu

Synthetic proxy OOD data

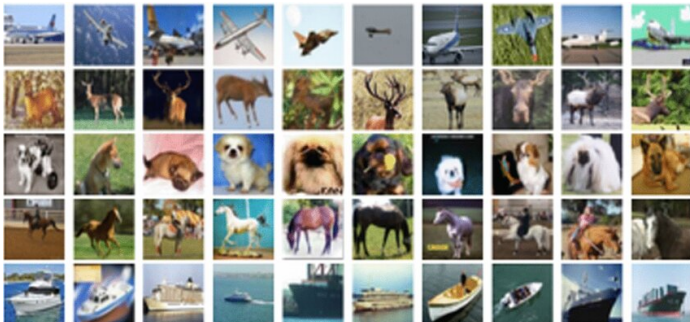
CSI: Novelty Detection via Contrastive Learning on Distributionally Shifted Instances

Jihoon Tack^{*†}, Sangwoo Mo^{*†}, Jongheon Jeong[‡], Jinwoo Shin^{†‡}

[†]Graduate School of AI, KAIST

[‡]School of Electrical Engineering, KAIST

Known outliers: TinyImages dataset (superset of CIFAR10/100)



Known outliers: synthetic image transformations



(a) Original

(b) Cutout

(c) Sobel

(d) Noise

(e) Blur

(f) Perm

(g) Rotate

Using proxy OOD data

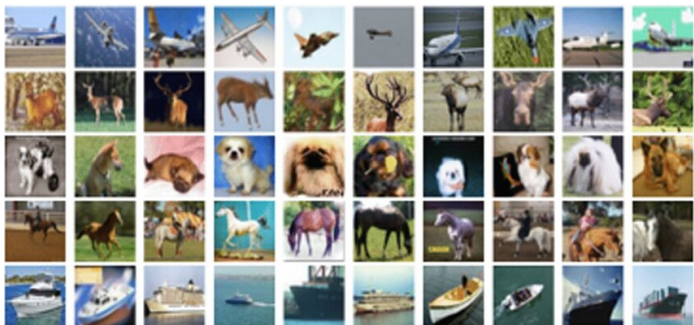
DEEP ANOMALY DETECTION WITH OUTLIER EXPOSURE

Dan Hendrycks
University of California, Berkeley
hendrycks@berkeley.edu

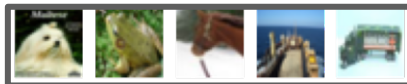
Mantas Mazeika
University of Chicago
mantas@ttic.edu

Thomas Dietterich
Oregon State University
tgd@oregonstate.edu

Known outliers: TinyImages dataset
(superset of CIFAR10/100)



In-distribution data:
5-class CIFAR10

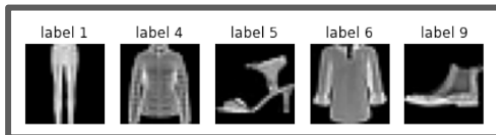


Outliers: remaining
CIFAR10 classes

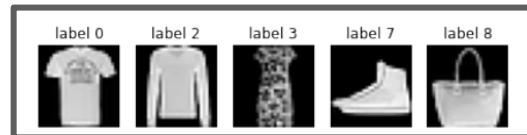


AUROC
Outlier exposure method: **0.82**

In-distribution data:
5-class FashionMNIST



Outliers: remaining
FashionMNIST classes



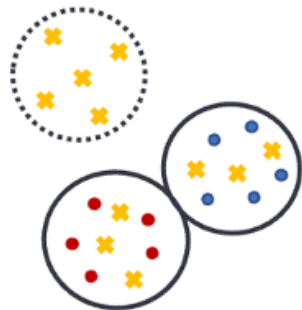
AUROC
Outlier exposure method: **0.66**

Semi-supervised OOD detection

Leveraging unlabeled data

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland



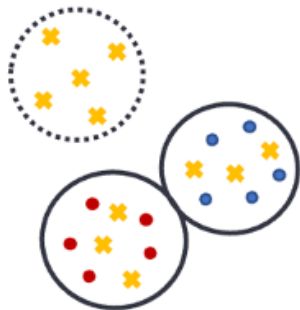
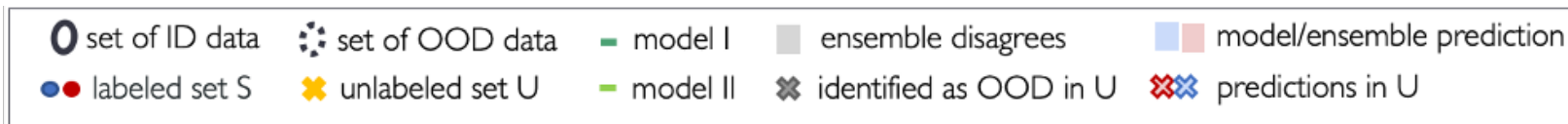
SSND setting

Semi-supervised OOD detection

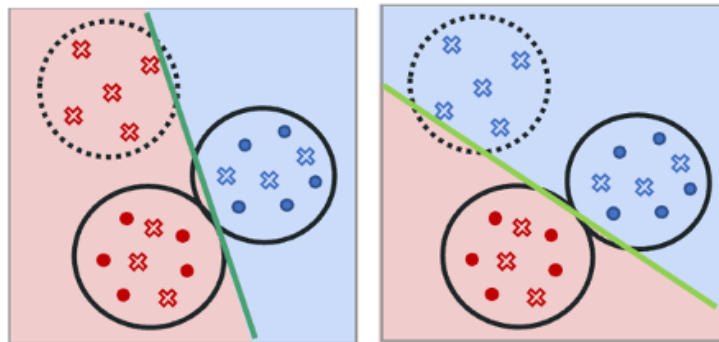
Leveraging unlabeled data

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
 Department of Computer Science
 ETH Zurich, Switzerland



SSND setting



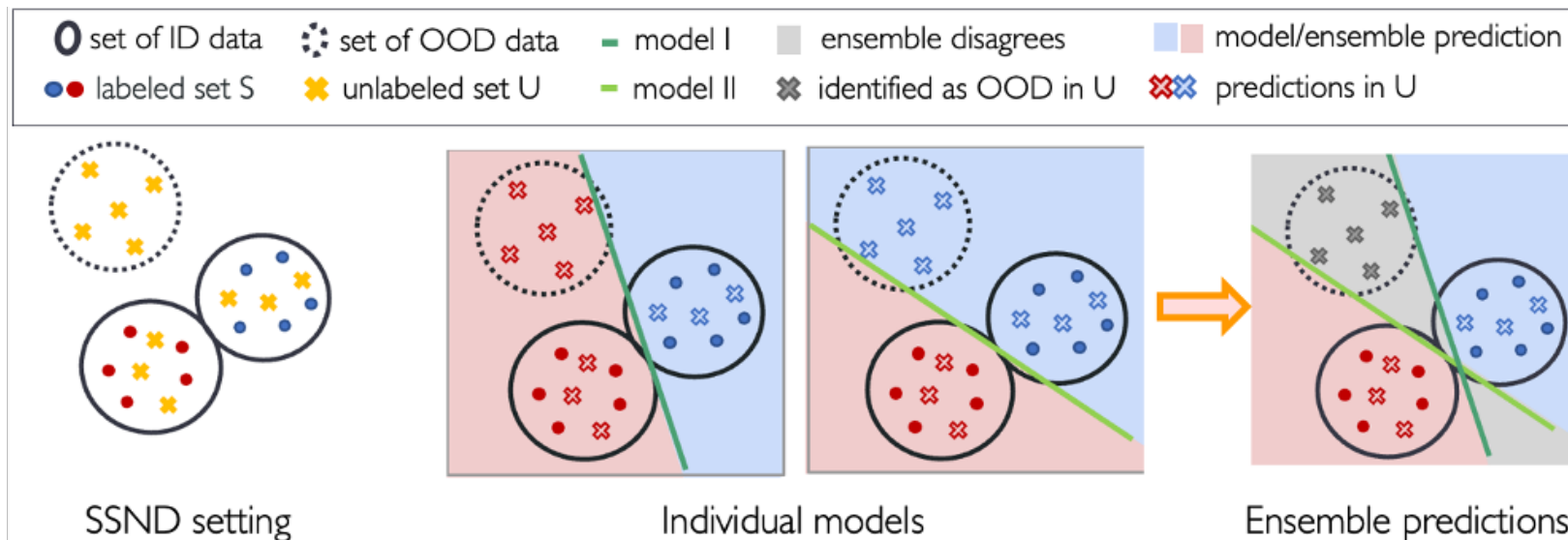
Individual models

Semi-supervised OOD detection

Leveraging unlabeled data

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
 Department of Computer Science
 ETH Zurich, Switzerland



sample x is flagged as OOD if “disagreement” $>$ threshold

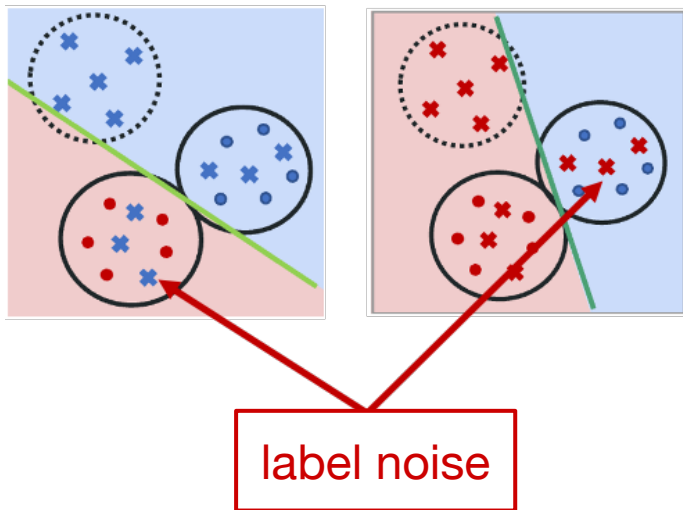
e.g. average pairwise TV distance between predictive distributions of the models in ensemble

Semi-supervised OOD detection

Key ingredient: Appropriate regularization

Semi-supervised novelty detection using
ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland

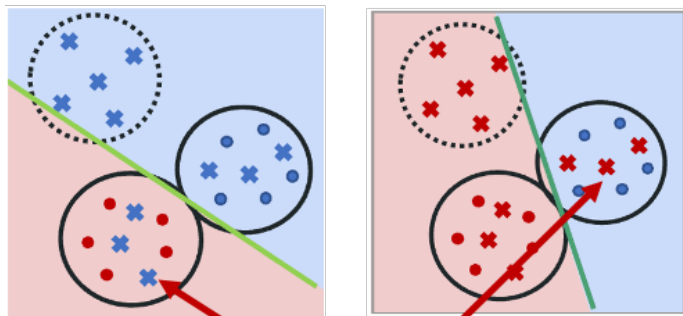


Semi-supervised OOD detection

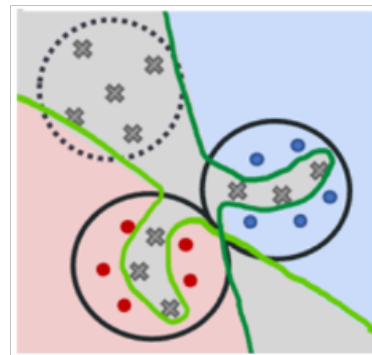
Key ingredient: Appropriate regularization

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland



label noise



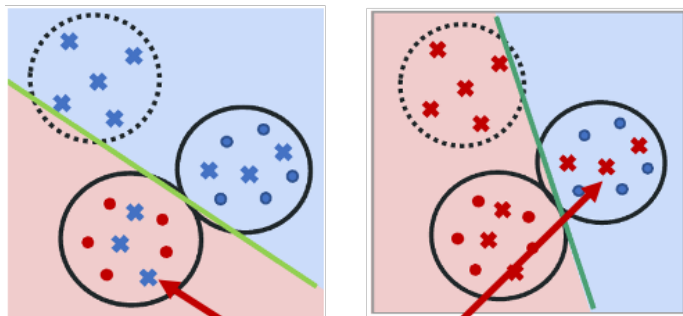
Too much diversity

Semi-supervised OOD detection

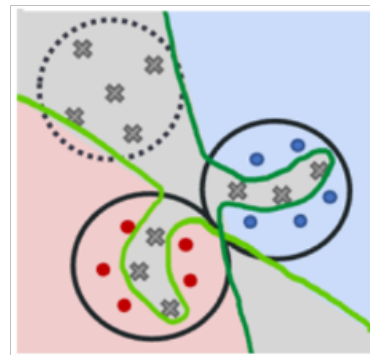
Key ingredient: Appropriate regularization

Semi-supervised novelty detection using ensembles with regularized disagreement

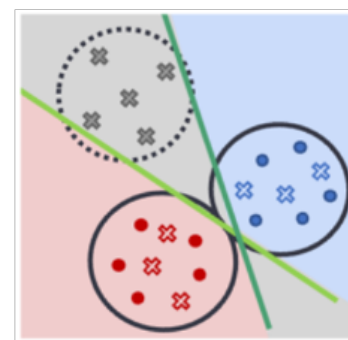
Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland



label noise



Too much diversity



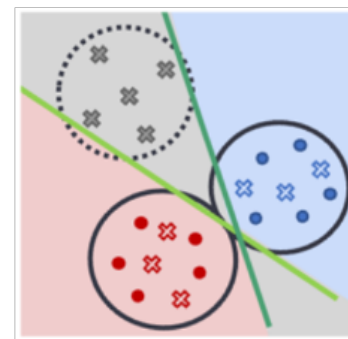
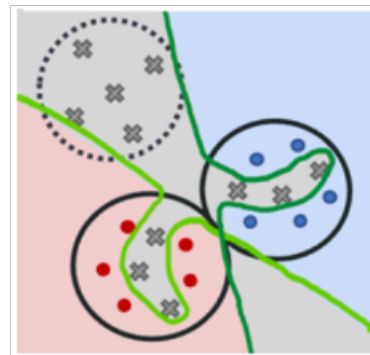
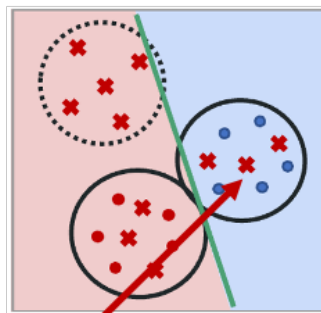
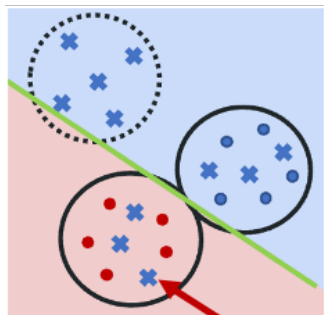
Right amount of diversity

Semi-supervised OOD detection

Key ingredient: Appropriate regularization

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țîfrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland



label noise

Too much diversity

Right amount of diversity

Idea: regularization with strength chosen using ID validation set

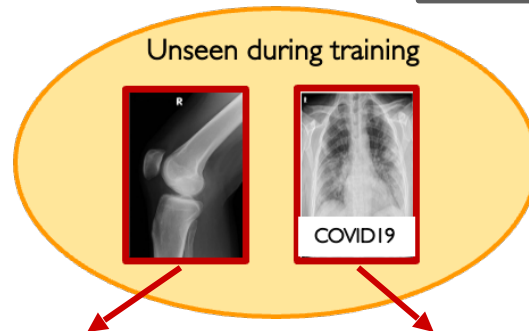
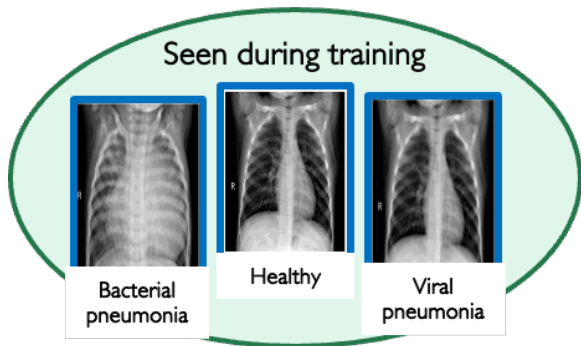
i.e. control FPR (ID samples incorrectly flagged as OOD)

Semi-supervised OOD detection

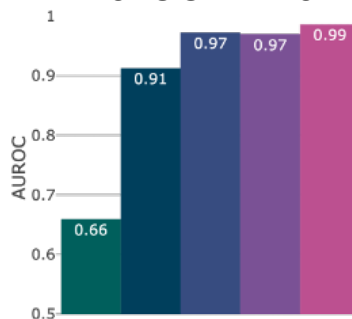
Performance on near OOD data

Semi-supervised novelty detection using ensembles with regularized disagreement

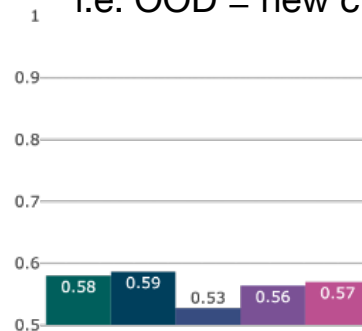
Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland



Far OOD (“easy”)
i.e. OOD = new dataset



Near OOD (“hard”)
i.e. OOD = new classes

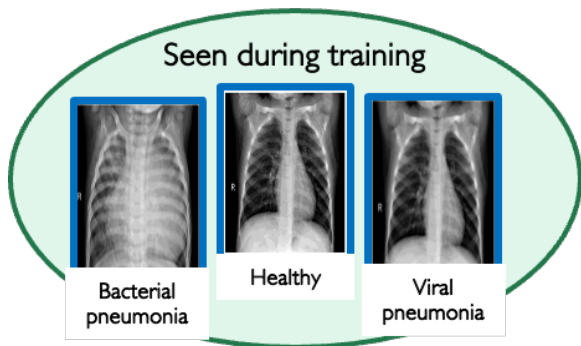


Semi-supervised OOD detection

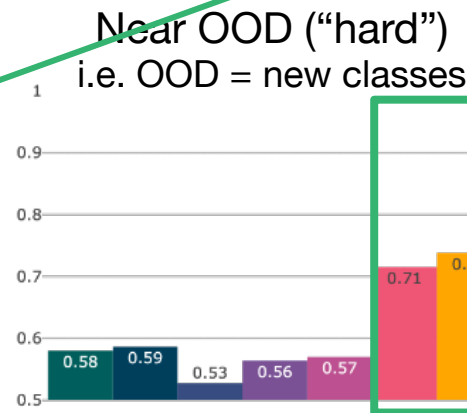
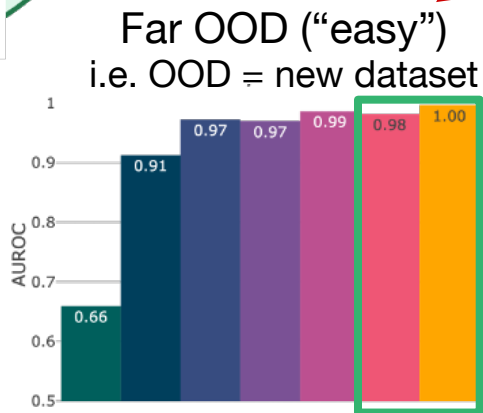
Performance on near OOD data

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland



semi-sup. OOD detection



Limitations of semi-supervised OOD detection

Challenge #1: not suitable for real-time applications

Semi-supervised novelty detection using
ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland

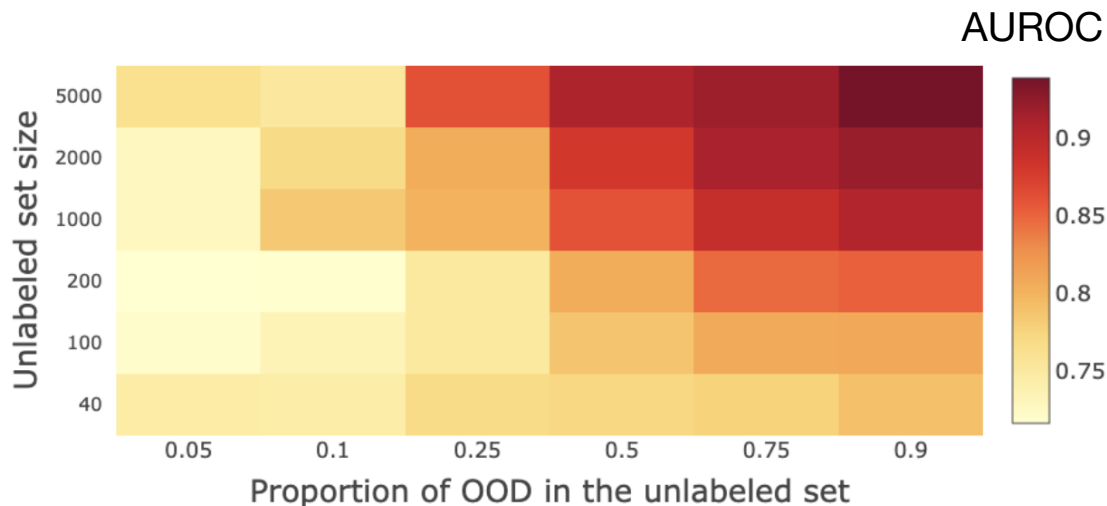
Limitations of semi-supervised OOD detection

Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang
Department of Computer Science
ETH Zurich, Switzerland

Challenge #1: not suitable for real-time applications

Challenge #2: not suitable for anomaly detection i.e. singleton outliers



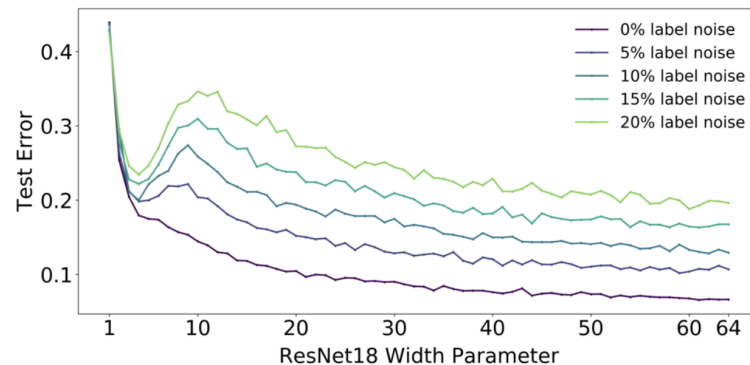
Outlook and future directions

Summary: Trustworthy ML under **imperfect data**

Summary: Trustworthy ML under imperfect data

If accuracy alone is the goal:

benign overfitting of label noise



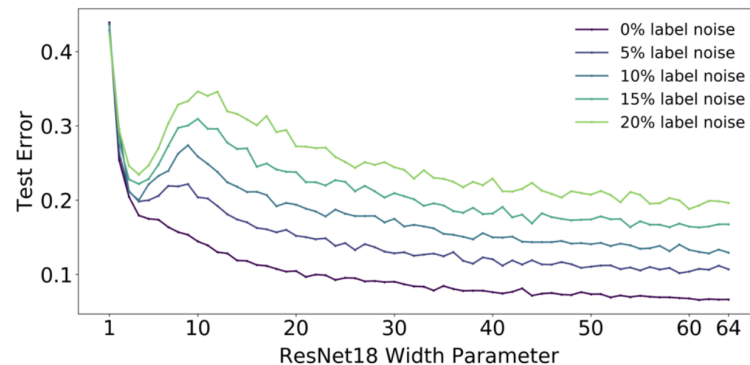
Summary: Trustworthy ML under imperfect data

If accuracy alone is the goal:

benign overfitting of label noise

If we care about trustworthiness:

This tutorial: Several examples of trustworthy learning algorithms that work well under label noise, missing data etc.



Summary: Trustworthy ML under imperfect data

If accuracy alone is the goal:

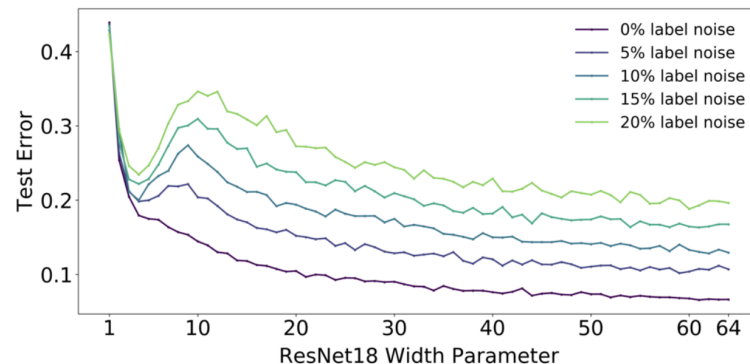
benign overfitting of label noise

If we care about trustworthiness:

This tutorial: Several examples of trustworthy learning algorithms that work well under label noise, missing data etc.

Open questions

- What other data-related limitations do existing trustworthy algorithms suffer from?
- How to improve trustworthiness in other difficult problem settings?



Summary: Trustworthy ML and **unlabeled data**

Summary: Trustworthy ML and **unlabeled data**

If accuracy alone is the goal:

SSL cannot be simultaneously better than both unsupervised and supervised learning

Can semi-supervised learning use all the data effectively? A lower bound perspective

Alexandru Tifrea*
ETH Zurich
alexandru.tifrea@inf.ethz.ch

Gizem Yüce*
EPFL
gizem.yuce@epfl.ch

Amartya Sanyal
Max Planck Institute for Intelligent Systems, Tübingen
amsa@di.ku.dk

Fanny Yang
ETH Zurich
fan.yan@inf.ethz.ch

Summary: Trustworthy ML and **unlabeled data**

If accuracy alone is the goal:

SSL cannot be simultaneously better than both unsupervised and supervised learning

If we care about trustworthiness:

This tutorial: Several examples where unlabeled data can help to overcome limitations of supervised learning.

Can semi-supervised learning use all the data effectively? A lower bound perspective

Alexandru Tifrea*
ETH Zurich
alexandru.tifrea@inf.ethz.ch

Gizem Yüce*
EPFL
gizem.yuce@epfl.ch

Amartya Sanyal
Max Planck Institute for Intelligent Systems, Tübingen
amsa@di.ku.dk

Fanny Yang
ETH Zurich
fan.yan@inf.ethz.ch

Summary: Trustworthy ML and **unlabeled data**

If accuracy alone is the goal:

SSL cannot be simultaneously better than both unsupervised and supervised learning

If we care about trustworthiness:

This tutorial: Several examples where unlabeled data can help to overcome limitations of supervised learning.

Open questions

- How fundamental are the improvements to trustworthiness due to unlabeled data?
- What other kinds of (potentially noisy) side information can be used to improve trustworthiness?

Can semi-supervised learning use all the data effectively? A lower bound perspective

Alexandru Tifrea*
ETH Zurich
alexandru.tifrea@inf.ethz.ch

Gizem Yüce*
EPFL
gizem.yuce@epfl.ch

Amartya Sanyal
Max Planck Institute for Intelligent Systems, Tübingen
amsa@di.ku.dk

Fanny Yang
ETH Zurich
fan.yan@inf.ethz.ch

Thank you!