

# **FR-Train: A Mutual Information-based Fair and Robust Training**

Yuji Roh, Kangwook Lee, Steven E. Whang, Changho Suh

**Yuji Roh**, Data Intelligence Lab, KAIST

# Trustworthy AI

“AI has significant potential to help solve challenging problems, including by advancing medicine, understanding language, and fueling scientific discovery. To realize that potential, it’s critical that AI is used and developed **responsibly**.”

-  AI, 2020

---

“Moving forward, “build for performance” will not suffice as an AI design paradigm. We must learn how to build, evaluate and monitor for **trust**.”

-  Trusting AI, 2020

# Trustworthy AI



**Fairness**



**Robustness**



**Explainability**



**Value  
Alignment**



**Transparency  
& Accountability**

# Trustworthy AI

Data-related



**Fairness**



**Robustness**



**Explainability**



**Value  
Alignment**



**Transparency  
& Accountability**

# Two approaches

- **Two-step** approach: Sanitize data -> Fair training  
Downside: very difficult to “decouple” poisoning and bias

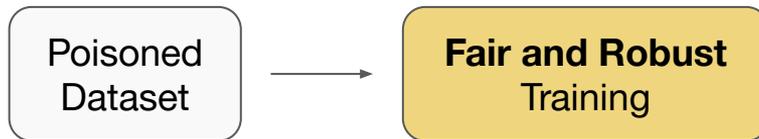


# Two approaches

- **Two-step** approach: Sanitize data -> Fair training  
Downside: very difficult to “decouple” poisoning and bias

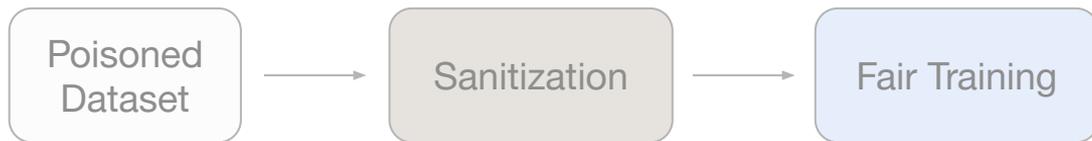


- **Holistic** approach: Fair & Robust training  
Performing the two operations along with model training results in much better performance



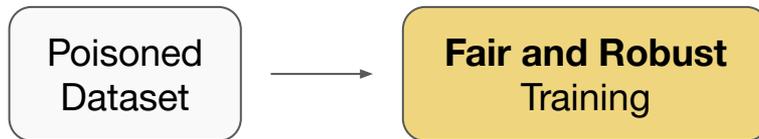
# Two approaches

- **Two-step** approach: Sanitize data -> Fair training  
Downside: very difficult to “decouple” poisoning and bias



**FR-Train**

- **Holistic** approach: Fair & Robust training  
Performing the two operations along with model training results in much better performance



**01**

Motivation

**02**

FR-Train

**03**

Experiments

**04**

Takeaways

**01**

**Motivation**

**02**

FR-Train

**03**

Experiments

**04**

Takeaways

# Trustworthy AI

Data-related



**Fairness**



**Robustness**



**Explainability**



**Value  
Alignment**



**Transparency  
& Accountability**

# Fairness

$X$	Feature
$Y$	Label
$Z$	Group attribute
$\hat{Y}$	Predicted label

- A machine learning model learns bias and discriminations in the data
- The fairness of a (binary) classifier can be defined in various ways:

## Demographic Parity

( $\Leftrightarrow$  Disparate Impact)

$$\mathbb{P}(\hat{Y} = 1|Z = 0) \approx \mathbb{P}(\hat{Y} = 1|Z = 1)$$

## Equalized Odds

$$\mathbb{P}(\hat{Y} = 1|Z = 0, Y = 1) \approx \mathbb{P}(\hat{Y} = 1|Z = 1, Y = 1)$$

$$\mathbb{P}(\hat{Y} = 1|Z = 0, Y = 0) \approx \mathbb{P}(\hat{Y} = 1|Z = 1, Y = 0)$$

- The level of fairness can be measured as a ratio or difference

# Fairness

$X$	Feature
$Y$	Label
$Z$	Group attribute
$\hat{Y}$	Predicted label

- A machine learning model learns bias and discriminations in the data
- The fairness of a (binary) classifier can be defined in various ways:

## Demographic Parity

( $\Leftrightarrow$  Disparate Impact)

$$\mathbb{P}(\hat{Y} = 1|Z = 0) \approx \mathbb{P}(\hat{Y} = 1|Z = 1)$$

## Equalized Odds

$$\mathbb{P}(\hat{Y} = 1|Z = 0, Y = 1) \approx \mathbb{P}(\hat{Y} = 1|Z = 1, Y = 1)$$

$$\mathbb{P}(\hat{Y} = 1|Z = 0, Y = 0) \approx \mathbb{P}(\hat{Y} = 1|Z = 1, Y = 0)$$

 **In this talk** 

- The level of fairness can be measured as a ratio or difference

$$DI := \min\left(\frac{\mathbb{P}(\hat{Y}=1|Z=0)}{\mathbb{P}(\hat{Y}=1|Z=1)}, \frac{\mathbb{P}(\hat{Y}=1|Z=1)}{\mathbb{P}(\hat{Y}=1|Z=0)}\right)$$

# Robustness

- Datasets are **easy to publish** nowadays, but as a result **easy to “poison”** as well
  - Poison = noisy, subjective, or even adversarial

kaggle

Google  
Dataset Search

- Attacker's goal : Increase the test loss by poisoning data
- Defender's goal : Train a classifier with small test loss
- Already a serious issue in federated learning

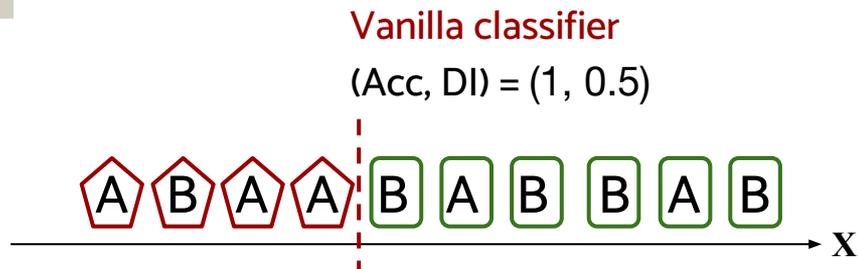
# Fairness + Robustness

What happens if we just apply a fairness-aware algorithm on a poisoned dataset?

- May result in a strictly **suboptimal** (accuracy, fairness) than vanilla training

# Motivating example

Clean



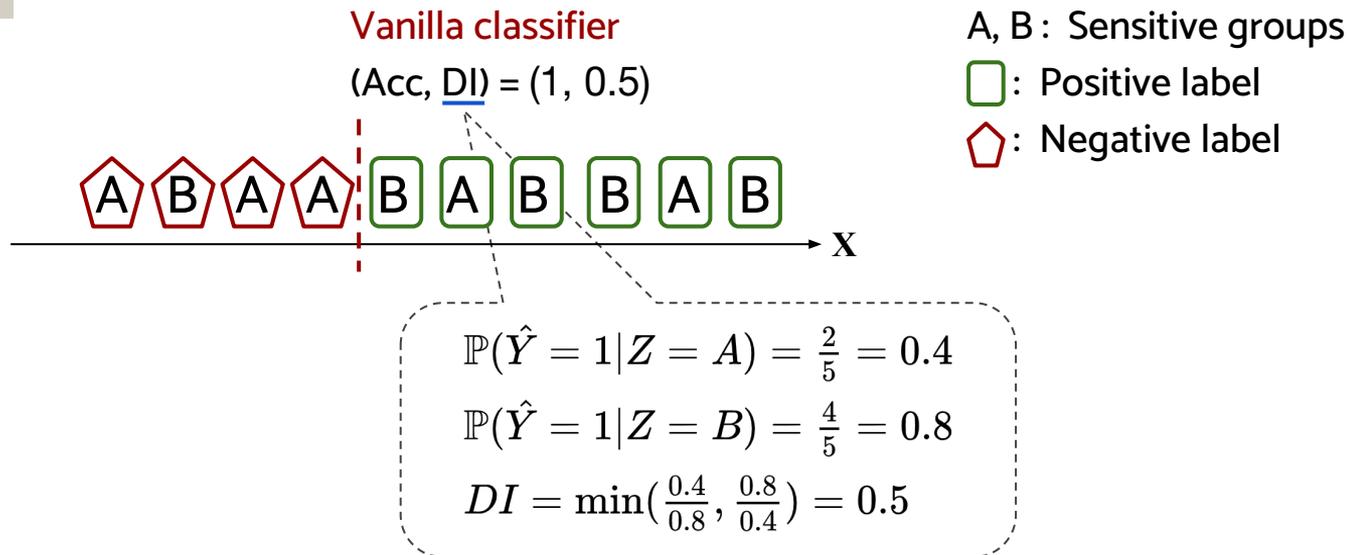
A, B: Sensitive groups

□: Positive label

⬠: Negative label

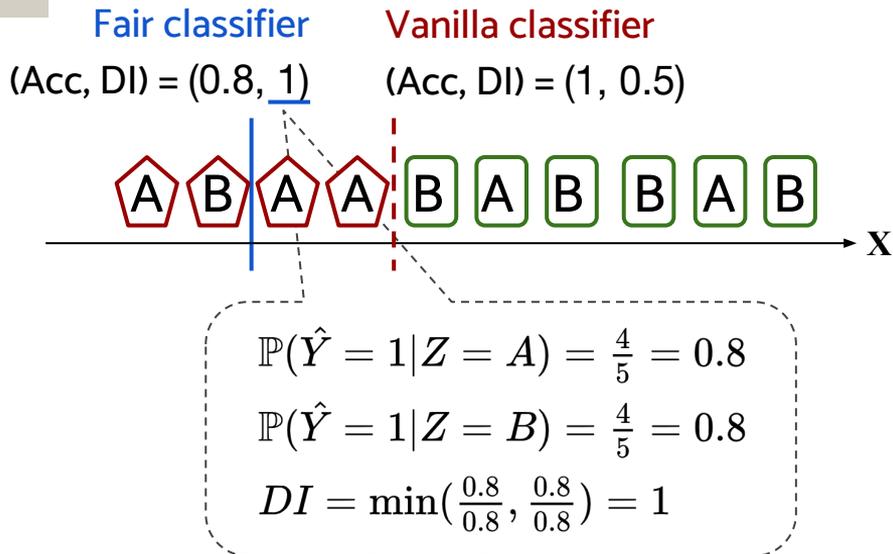
# Motivating example

Clean



# Motivating example

Clean



A, B: Sensitive groups

□: Positive label

⬠: Negative label

# Motivating example

Clean

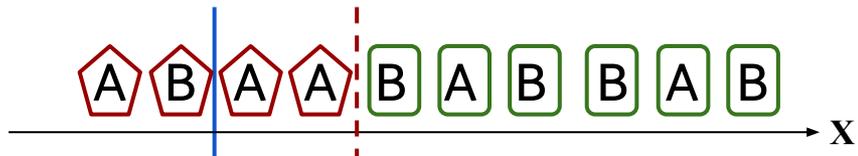
Fair classifier  
(Acc, DI) = (0.8, 1)

Vanilla classifier  
(Acc, DI) = (1, 0.5)

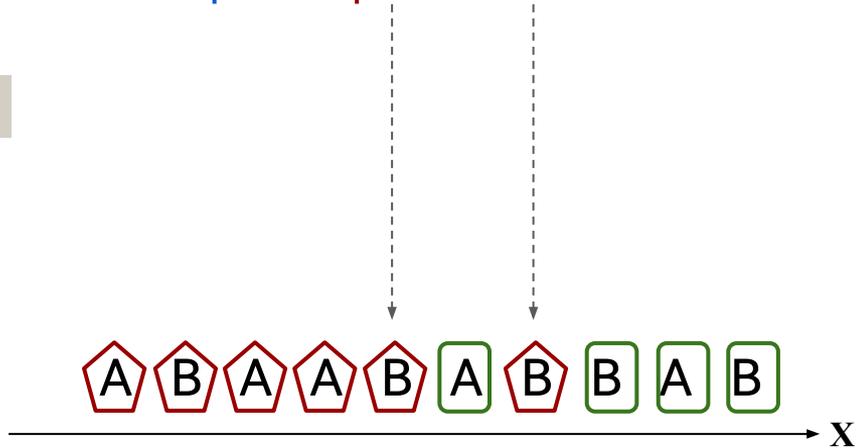
A, B: Sensitive groups

□: Positive label

⬠: Negative label

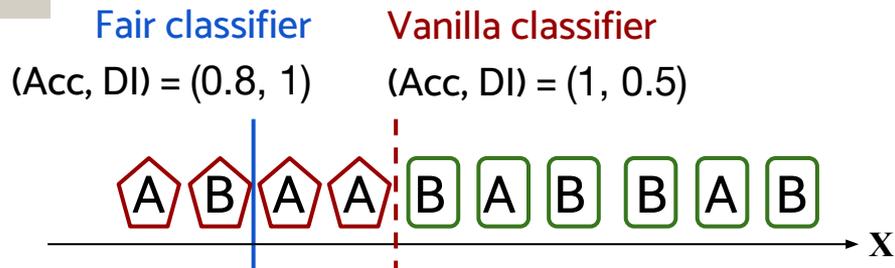


Poisoned

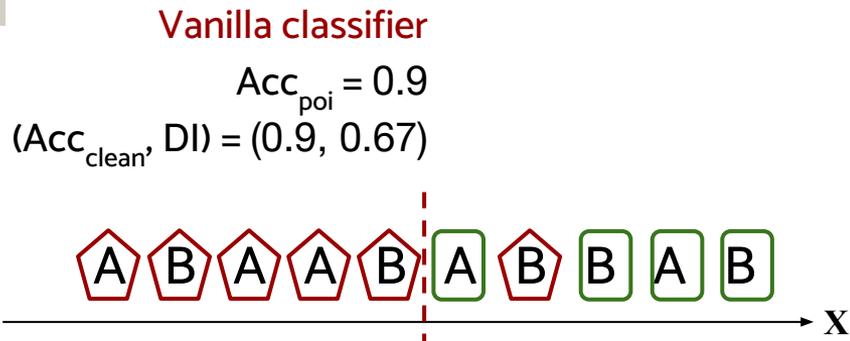


# Motivating example

Clean



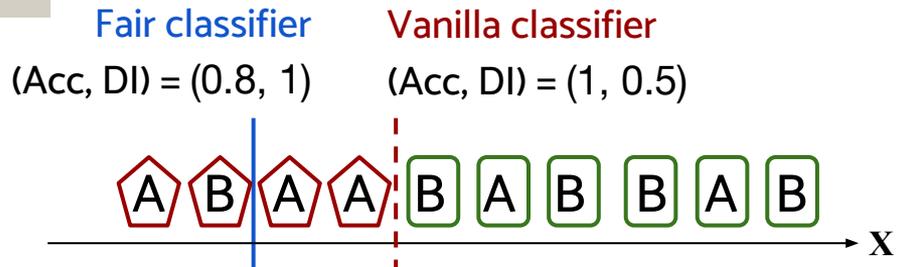
Poisoned



Acc: ↓  
DI: ↑

# Motivating example

Clean

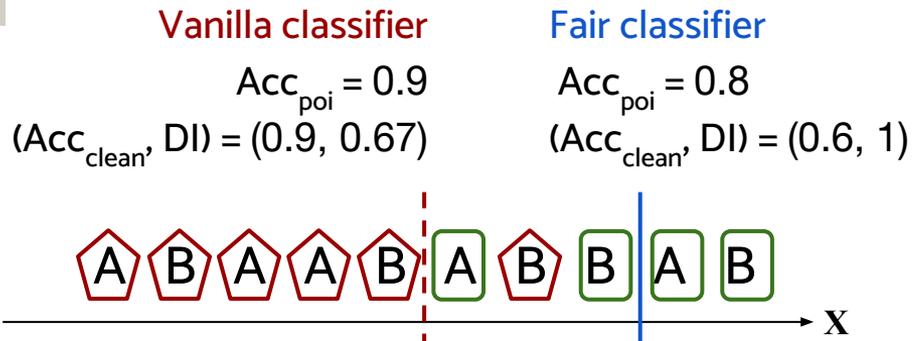


A, B: Sensitive groups

□: Positive label

⬠: Negative label

Poisoned



Acc: ↓  
DI: —

Suboptimal

# Fairness + Robustness

What happens if we just apply a fairness-aware algorithm on a poisoned dataset?

- May result in a strictly **suboptimal** (accuracy, fairness) than vanilla training

We need a holistic approach to fair and robust training. **FR-Train!**

01

Motivation

02

**FR-Train**

03

Experiments

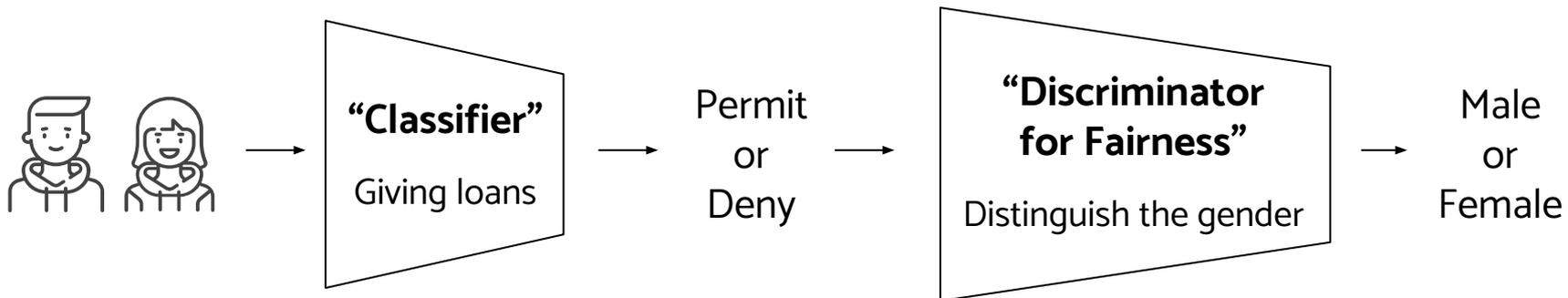
04

Takeaways

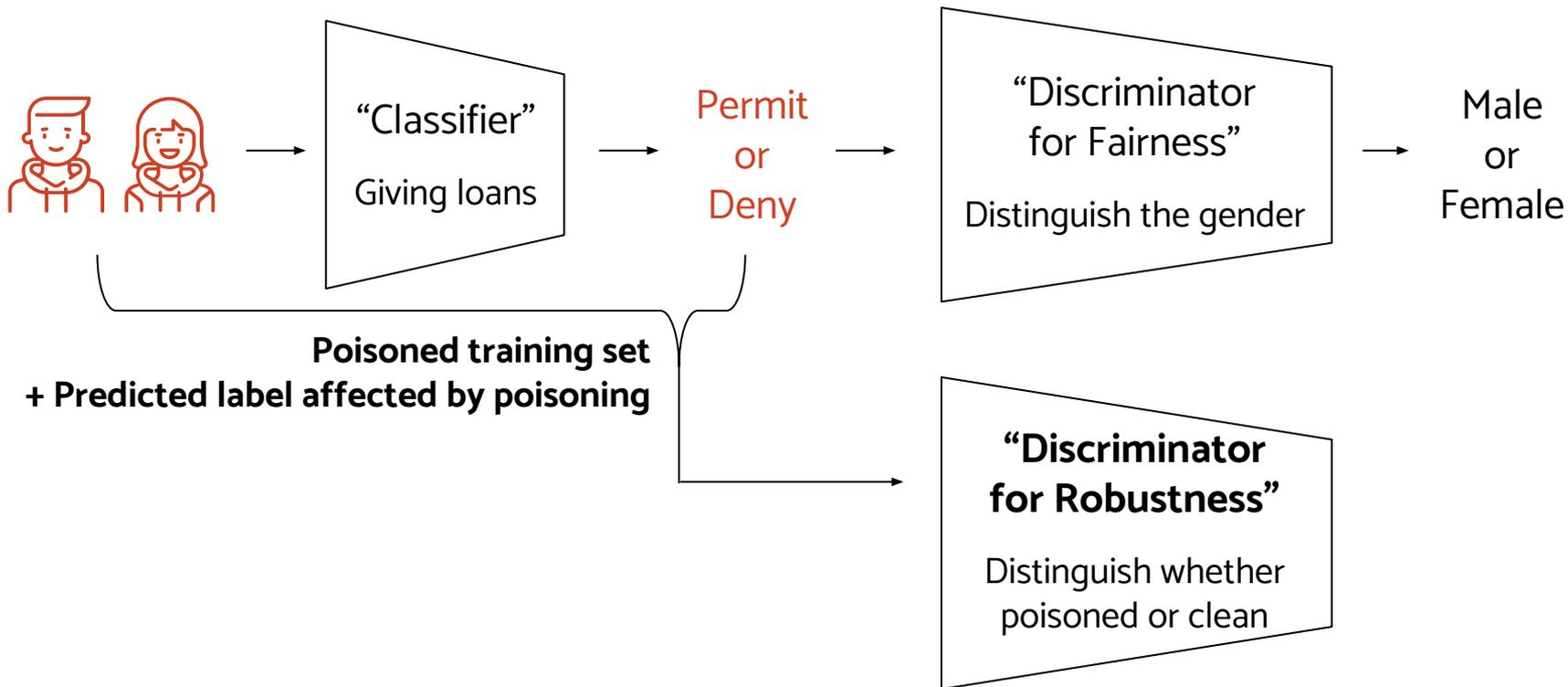
# FR-Train - Main contributions

- FR-Train is a **holistic framework** for fair and robust training
- Extends a state-of-the-art fairness-only method called Adversarial Debiasing
  - Provides a novel mutual information (MI)-based interpretation of adversarial learning
  - Adds a robust discriminator that uses a small clean validation set for data sanitization
- We also propose crowdsourcing methods for constructing a clean validation set

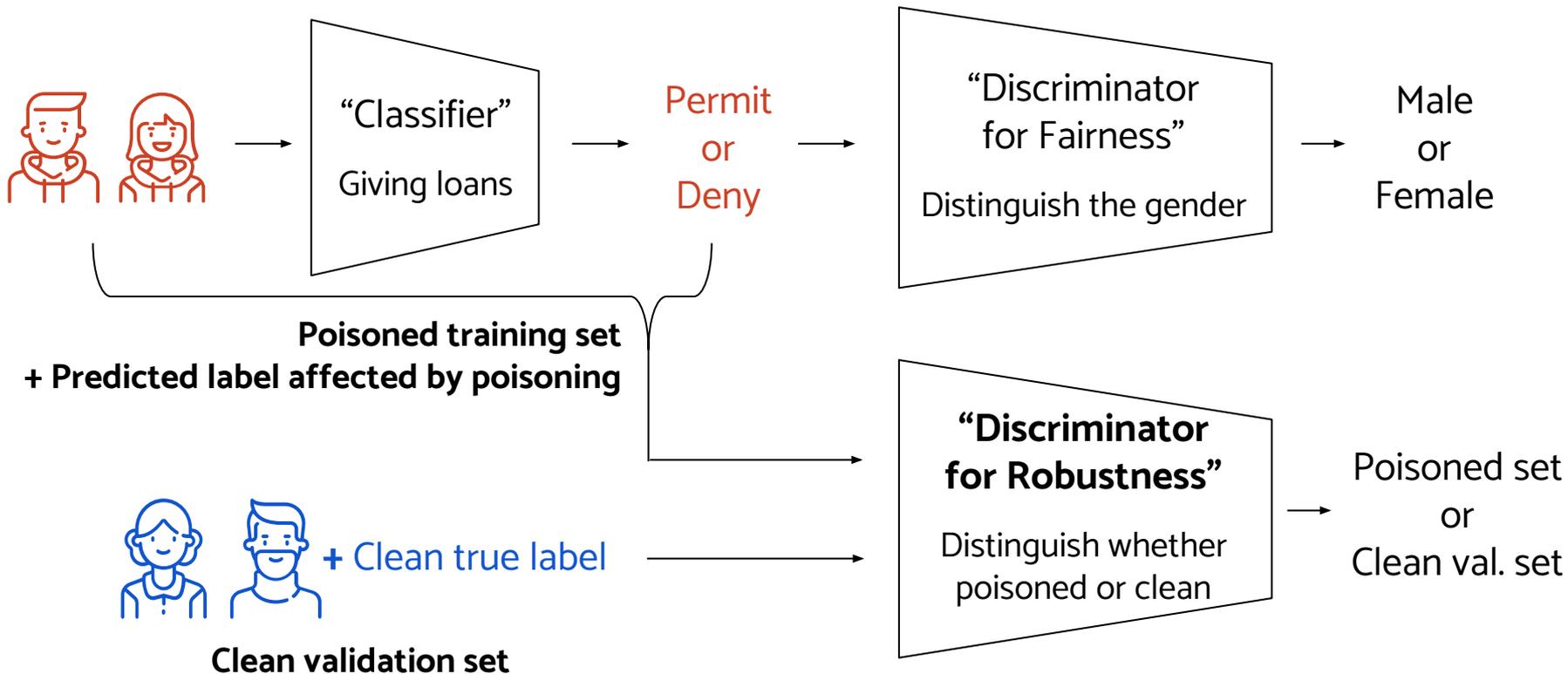
# FR-Train



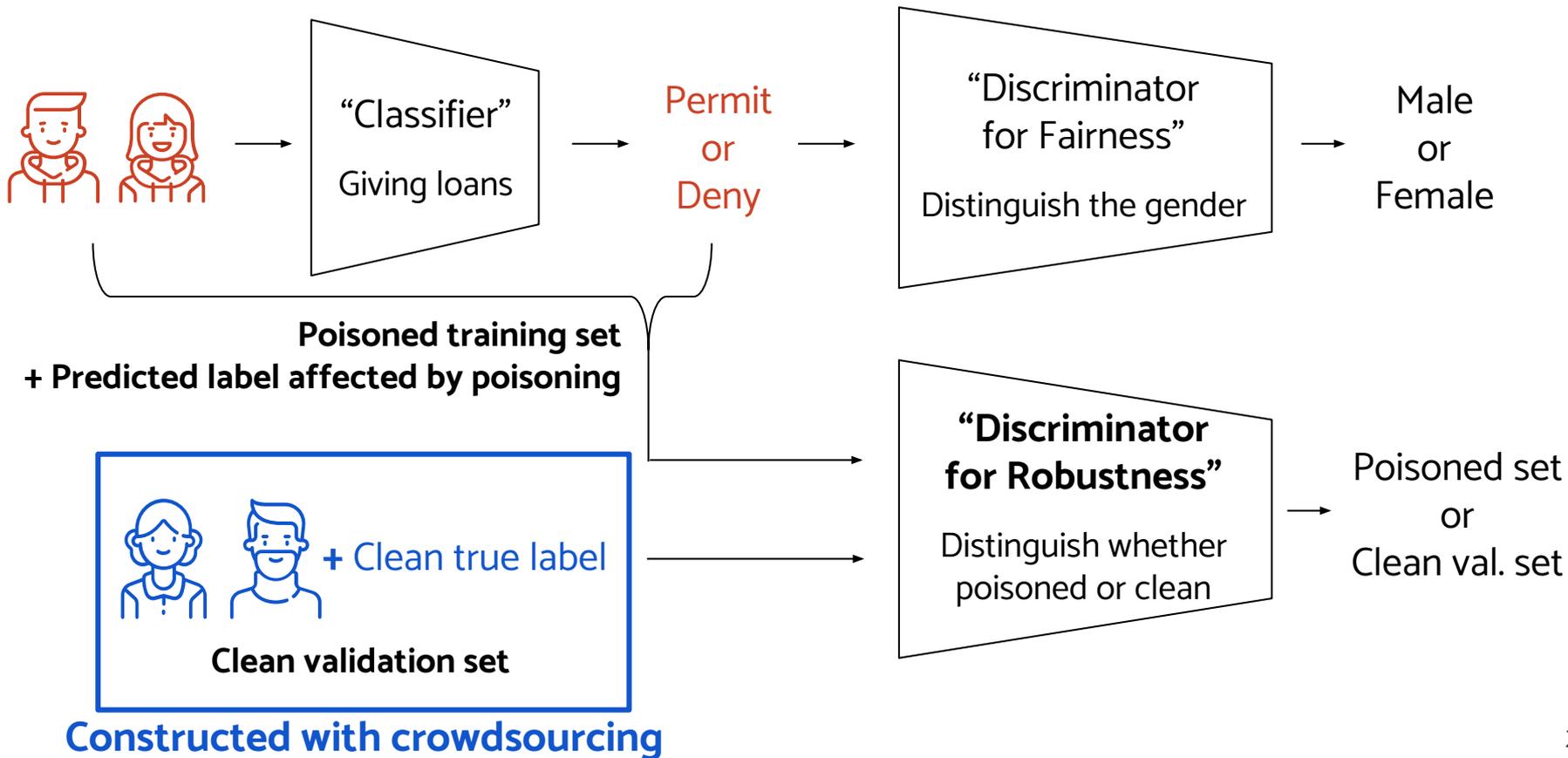
# FR-Train



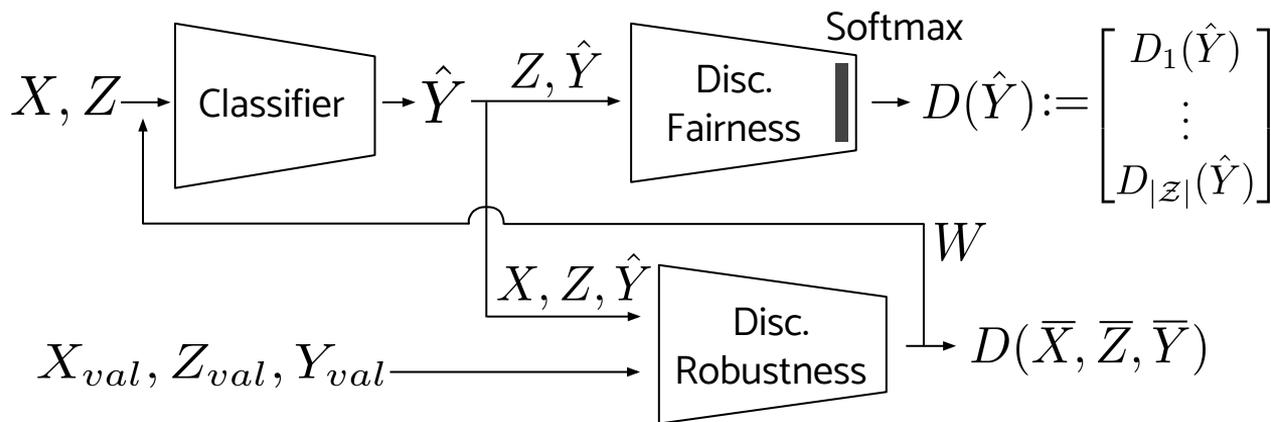
# FR-Train



# FR-Train



# Mutual information-based interpretation



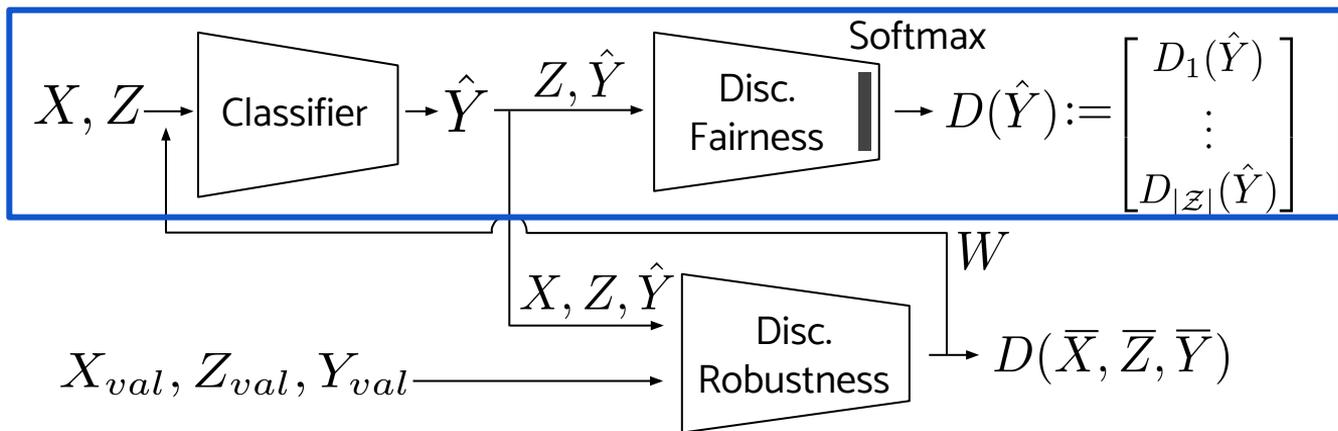
## Theorem 1 - Fairness

$$I(Z; \hat{Y}) = \max_{D_z(\hat{y}): \sum_z D_z(\hat{y})=1, \forall \hat{y}} \sum_{z \in \mathcal{Z}} P_Z(z) \mathbb{E}_{P_{\hat{Y}|z}} \left[ \log D_z(\hat{Y}) \right] + H(Z)$$

## Theorem 2 - Robustness

$$I(V; \bar{X}, \bar{Z}, \bar{Y}) = \max_{D_v(x,z,y): \sum_v D_v(x,z,y)=1, \forall (x,z,y)} \sum_{v \in \mathcal{V}} P_V(v) \mathbb{E}_{P_{\bar{X}, \bar{Z}, \bar{Y}|v}} \left[ \log D_v(\bar{X}, \bar{Z}, \bar{Y}) \right] + H(V)$$

# Mutual information-based interpretation



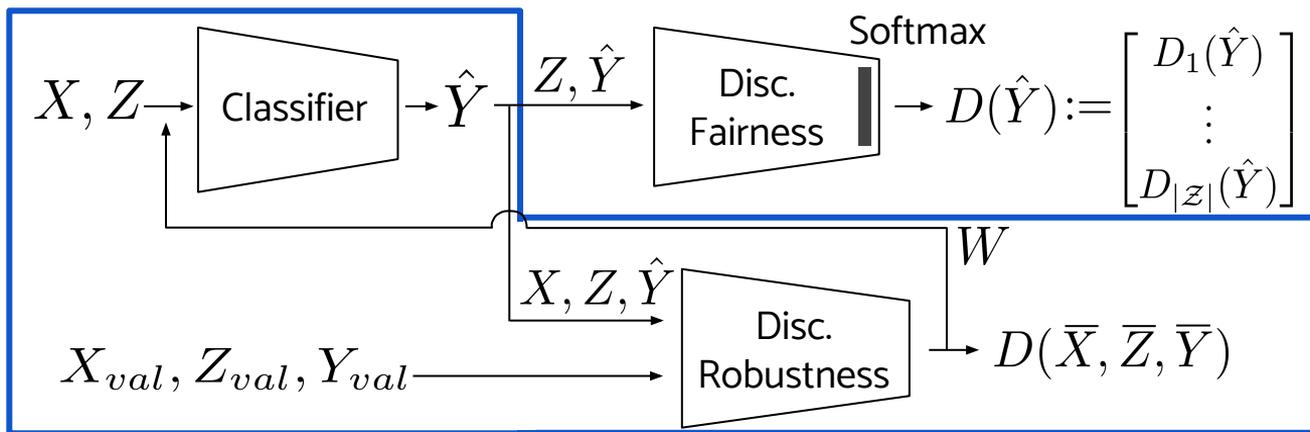
## Theorem 1 - Fairness

$$I(Z; \hat{Y}) = \max_{D_z(\hat{y}): \sum_z D_z(\hat{y})=1, \forall \hat{y}} \sum_{z \in \mathcal{Z}} P_Z(z) \mathbb{E}_{P_{\hat{Y}|z}} \left[ \log D_z(\hat{Y}) \right] + H(Z)$$

## Theorem 2 - Robustness

$$I(V; \bar{X}, \bar{Z}, \bar{Y}) = \max_{D_v(x,z,y): \sum_v D_v(x,z,y)=1, \forall (x,z,y)} \sum_{v \in \mathcal{V}} P_V(v) \mathbb{E}_{P_{\bar{X}, \bar{Z}, \bar{Y}|v}} \left[ \log D_v(\bar{X}, \bar{Z}, \bar{Y}) \right] + H(V)$$

# Mutual information-based interpretation



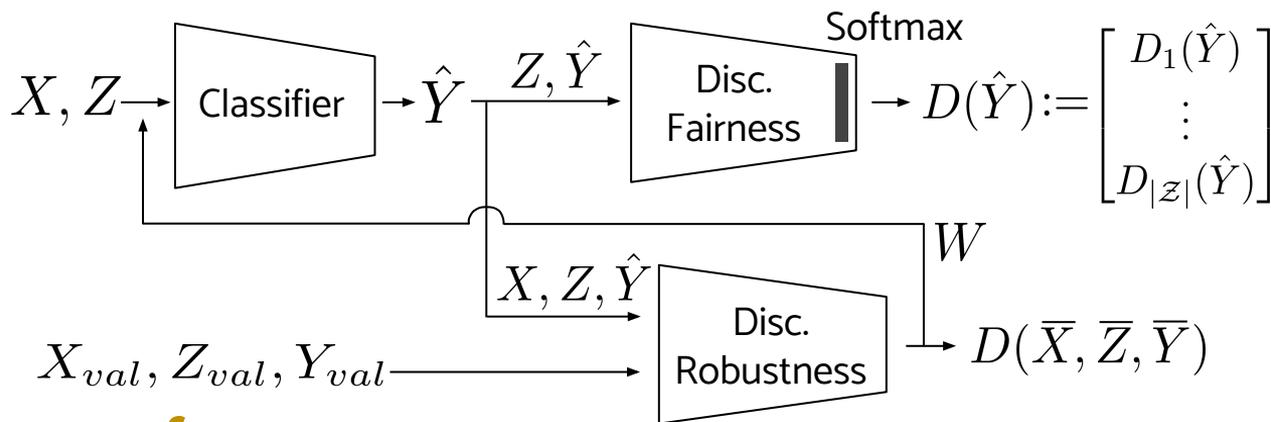
## Theorem 1 - Fairness

$$I(Z; \hat{Y}) = \max_{D_z(\hat{y}): \sum_z D_z(\hat{y})=1, \forall \hat{y}} \sum_{z \in \mathcal{Z}} P_Z(z) \mathbb{E}_{P_{\hat{Y}|z}} [\log D_z(\hat{Y})] + H(Z)$$

## Theorem 2 - Robustness

$$I(V; \bar{X}, \bar{Z}, \bar{Y}) = \max_{D_v(x,z,y): \sum_v D_v(x,z,y)=1, \forall (x,z,y)} \sum_{v \in \mathcal{V}} P_V(v) \mathbb{E}_{P_{\bar{X}, \bar{Z}, \bar{Y}|v}} [\log D_v(\bar{X}, \bar{Z}, \bar{Y})] + H(V)$$

# Mutual information-based interpretation



See paper for proofs

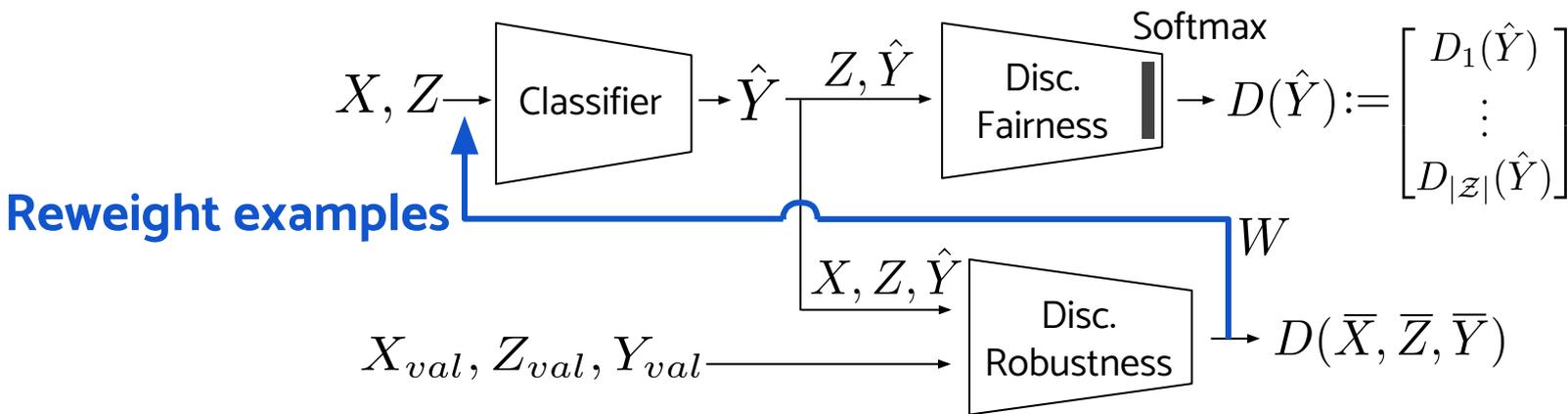
## Theorem 1 - Fairness

$$I(Z; \hat{Y}) = \max_{D_z(\hat{y}): \sum_z D_z(\hat{y})=1, \forall \hat{y}} \sum_{z \in \mathcal{Z}} P_Z(z) \mathbb{E}_{P_{\hat{Y}|z}} \left[ \log D_z(\hat{Y}) \right] + H(Z)$$

## Theorem 2 - Robustness

$$I(V; \bar{X}, \bar{Z}, \bar{Y}) = \max_{D_v(x,z,y): \sum_v D_v(x,z,y)=1, \forall (x,z,y)} \sum_{v \in \mathcal{V}} P_V(v) \mathbb{E}_{P_{\bar{X}, \bar{Z}, \bar{Y}|v}} \left[ \log D_v(\bar{X}, \bar{Z}, \bar{Y}) \right] + H(V)$$

# Mutual information-based interpretation



## Theorem 1 - Fairness

$$I(Z; \hat{Y}) = \max_{D_z(\hat{y}): \sum_z D_z(\hat{y})=1, \forall \hat{y}} \sum_{z \in \mathcal{Z}} P_Z(z) \mathbb{E}_{P_{\hat{Y}|z}} \left[ \log D_z(\hat{Y}) \right] + H(Z)$$

## Theorem 2 - Robustness

$$I(V; \bar{X}, \bar{Z}, \bar{Y}) = \max_{D_v(x,z,y): \sum_v D_v(x,z,y)=1, \forall (x,z,y)} \sum_{v \in \mathcal{V}} P_V(v) \mathbb{E}_{P_{\bar{X}, \bar{Z}, \bar{Y}|v}} \left[ \log D_v(\bar{X}, \bar{Z}, \bar{Y}) \right] + H(V)$$

01

Motivation

02

FR-Train

03

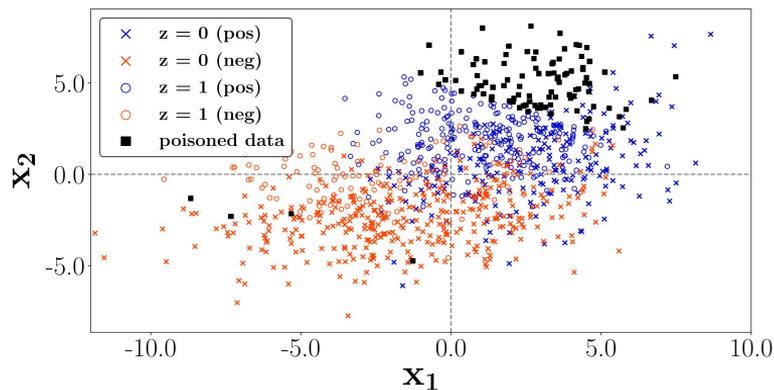
**Experiments**

04

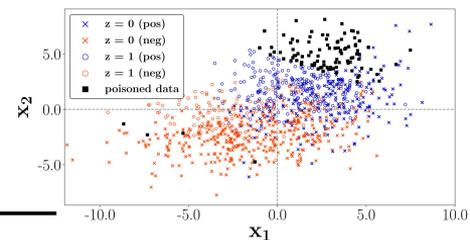
Takeaways

# Experimental setting

- Synthetic data
  - Poisoning (label flipping): 10% of training data
  - Validation set: 10% of training data
- Real data (results in paper)
  - **COMPAS**: Predict recidivism in two years for criminals
  - **AdultCensus**: Predict whether annual income > \$50K or not
  - Poisoning: 10% of training data
  - Validation set: 5% of training data

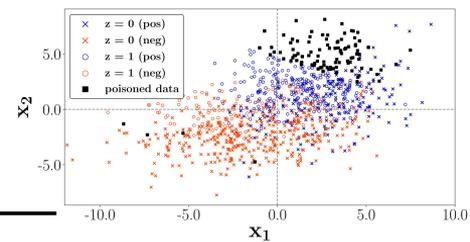


# Synthetic data results



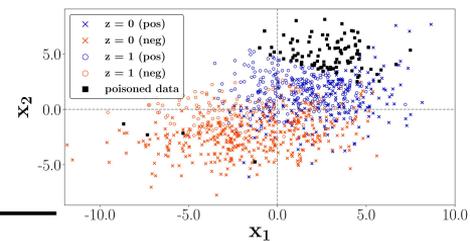
		Method	Clean data		Poisoned data	
			DI	Acc.	DI	Acc.
Fair-only algorithms	{	FC	.822	.806	.831 (1.1% ↑)	.760 (5.7% ↓)
		LBC	.819	.760	.827 (1.0% ↑)	.715 (5.9% ↓)
		AD	.807	.811	.834 (3.4% ↑)	.769 (5.2% ↓)
Two-step approach : Data sanitization + Fair training	{	RML+FC	.822	.806	.802 (2.4% ↓)	.529 (34.% ↓)
		RML+LBC	.819	.760	.810 (1.1% ↓)	.752 (1.1% ↓)
		RML+AD	.807	.811	.808 (0.1% ↑)	.756 (6.8% ↓)
Logistic regression	—	LR	.409	.885	.446 (9.1% ↑)	.819 (7.5% ↓)
Data sanitization	—	RML	.471	.876	.395 (16.% ↓)	.869 (0.8% ↓)
using clean val. set		<b>FR-Train</b>	<b>.818</b>	<b>.807</b>	<b>.827 (1.1% ↑)</b>	<b>.814 (0.9% ↑)</b>

# Synthetic data results



	Method	Clean data		Poisoned data		
		DI	Acc.	DI	Acc.	
Fair-only algorithms	FC	.822	.806	.831 (1.1% ↑)	.760 (5.7% ↓)	Low accuracy
	LBC	.819	.760	.827 (1.0% ↑)	.715 (5.9% ↓)	
	AD	.807	.811	.834 (3.4% ↑)	.769 (5.2% ↓)	
Two-step approach : Data sanitization + Fair training	RML+FC	.822	.806	.802 (2.4% ↓)	.529 (34.% ↓)	
	RML+LBC	.819	.760	.810 (1.1% ↓)	.752 (1.1% ↓)	
	RML+AD	.807	.811	.808 (0.1% ↑)	.756 (6.8% ↓)	
Logistic regression	LR	.409	.885	.446 (9.1% ↑)	.819 (7.5% ↓)	
Data sanitization	RML	.471	.876	.395 (16.% ↓)	.869 (0.8% ↓)	
using clean val. set	<b>FR-Train</b>	<b>.818</b>	<b>.807</b>	<b>.827 (1.1% ↑)</b>	<b>.814 (0.9% ↑)</b>	

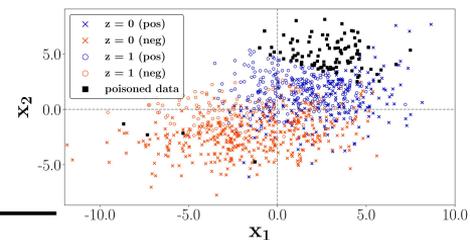
# Synthetic data results



	Method	Clean data		Poisoned data	
		DI	Acc.	DI	Acc.
Fair-only algorithms	FC	.822	.806	.831 (1.1% ↑)	.760 (5.7% ↓)
	LBC	.819	.760	.827 (1.0% ↑)	.715 (5.9% ↓)
	AD	.807	.811	.834 (3.4% ↑)	.769 (5.2% ↓)
Two-step approach : Data sanitization + Fair training	RML+FC	.822	.806	.802 (2.4% ↓)	.529 (34.% ↓)
	RML+LBC	.819	.760	.810 (1.1% ↓)	.752 (1.1% ↓)
	RML+AD	.807	.811	.808 (0.1% ↑)	.756 (6.8% ↓)
<b>Logistic regression</b>	LR	.409	.885	.446 (9.1% ↑)	.819 (7.5% ↓)
<b>Data sanitization</b>	RML	.471	.876	.395 (16.% ↓)	.869 (0.8% ↓)
<b>using clean val. set</b>	<b>FR-Train</b>	<b>.818</b>	<b>.807</b>	<b>.827 (1.1% ↑)</b>	<b>.814 (0.9% ↑)</b>

Low fairness

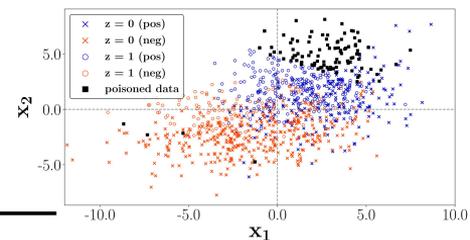
# Synthetic data results



		Method	Clean data		Poisoned data	
			DI	Acc.	DI	Acc.
Fair-only algorithms	{	FC	.822	.806	.831 (1.1% ↑)	.760 (5.7% ↓)
		LBC	.819	.760	.827 (1.0% ↑)	.715 (5.9% ↓)
		AD	.807	.811	.834 (3.4% ↑)	.769 (5.2% ↓)
<b>Two-step approach :</b> Data sanitization + Fair training	{	RML+FC	.822	.806	.802 (2.4% ↓)	<b>.529 (34.% ↓)</b>
		RML+LBC	.819	.760	.810 (1.1% ↓)	<b>.752 (1.1% ↓)</b>
		RML+AD	.807	.811	.808 (0.1% ↑)	<b>.756 (6.8% ↓)</b>
Logistic regression	—	LR	.409	.885	.446 (9.1% ↑)	.819 (7.5% ↓)
Data sanitization	—	RML	.471	.876	.395 (16.% ↓)	.869 (0.8% ↓)
using clean val. set		<b>FR-Train</b>	<b>.818</b>	<b>.807</b>	<b>.827 (1.1% ↑)</b>	<b>.814 (0.9% ↑)</b>

Also low accuracy

# Synthetic data results



		Method	Clean data		Poisoned data	
			DI	Acc.	DI	Acc.
Fair-only algorithms	{	FC	.822	.806	.831 (1.1% ↑)	.760 (5.7% ↓)
		LBC	.819	.760	.827 (1.0% ↑)	.715 (5.9% ↓)
		AD	.807	.811	.834 (3.4% ↑)	.769 (5.2% ↓)
Two-step approach : Data sanitization + Fair training	{	RML+FC	.822	.806	.802 (2.4% ↓)	.529 (34.% ↓)
		RML+LBC	.819	.760	.810 (1.1% ↓)	.752 (1.1% ↓)
		RML+AD	.807	.811	.808 (0.1% ↑)	.756 (6.8% ↓)
Logistic regression	—	LR	.409	.885	.446 (9.1% ↑)	.819 (7.5% ↓)
Data sanitization	—	RML	.471	.876	.395 (16.% ↓)	.869 (0.8% ↓)
using clean val. set		<b>FR-Train</b>	<b>.818</b>	<b>.807</b>	<b>.827 (1.1% ↑)</b>	<b>.814 (0.9% ↑)</b>

**Holistic approach =  
high fairness & accuracy**

01

Motivation

02

FR-Train

03

Experiments

04

**Takeaways**

# Takeaways

- Trustworthy AI needs both fair and robust training
- However, addressing fairness and robustness separately is suboptimal
- FR-Train is a **holistic framework for trustworthy AI** performing fair and robust training
  - Mutual information-based interpretation of adversarial learning
  - Novel architecture that enjoys the synergistic effect of fair and robust discriminators
  - Requires a small clean validation set, which can be constructed using crowdsourcing
- Lots of open problems
  - Without clean validation set
  - Other poisoning
  - Algorithm stability

# Takeaways

- Trustworthy AI needs both fair and robust training
- However, addressing fairness and robustness separately is suboptimal
- FR-Train is a **holistic framework for trustworthy AI** performing fair and robust training
  - Mutual information-based interpretation of adversarial learning
  - Novel architecture that enjoys the synergistic effect of fair and robust discriminators
  - Requires a small clean validation set, which can be constructed using crowdsourcing
- Lots of open problems
  - Without clean validation set
  - Other poisoning
  - Algorithm stability

Thank you :)