

Robust Learning from Multiple Sources

Christoph H. Lampert



Mathematical Machine Learning Seminar MPI Mis + UCLA
April 14th, 2022

Institute of Science and Technology Austria (ISTA)



- ▶ public research institute, opened in 2009
- ▶ located in outskirts of Vienna

Focus on curiosity-driven basic research

- ▶ avoiding boundaries between disciplines
- ▶ current 70 research groups
 - ▶ Computer Science, Mathematics, Physics, Chemistry, Biology, Neuroscience, Earth and Climate Sciences
- ▶ ELLIS unit since 2019

We're hiring! (on all levels)

- ▶ interns, PhD students, postdocs
- ▶ faculty (tenure-track or tenured), ...

More information: chl@ist.ac.at or <https://cvml.ist.ac.at>

Topics in Our Research Group

Machine Learning Theory

- ▶ Transfer Learning
- ▶ Multi-task Learning
- ▶ Lifelong/Meta-Learning
- ▶ Multi-source/Federated Learning

Models/Algorithms

- ▶ Zero-shot Learning
- ▶ Continual Learning
- ▶ Weakly-supervised Learning
- ▶ Trustworthy/Robust Learning

Learning for Computer Vision

- ▶ Scene Understanding
- ▶ Generative Models
- ▶ Abstract Reasoning
- ▶ Semantic Representations

Topics in Our Research Group

Machine Learning Theory

- ▶ Transfer Learning
- ▶ Multi-task Learning
- ▶ Lifelong/Meta-Learning
- ▶ Multi-source/Federated Learning

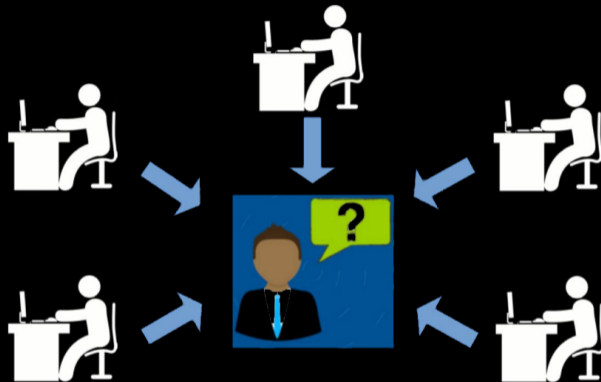
Models/Algorithms

- ▶ Zero-shot Learning
- ▶ Continual Learning
- ▶ Weakly-supervised Learning
- ▶ Trustworthy/Robust Learning

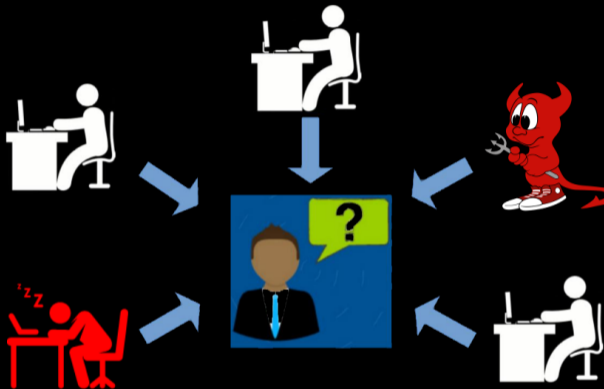
Learning for Computer Vision

- ▶ Scene Understanding
- ▶ Generative Models
- ▶ Abstract Reasoning
- ▶ Semantic Representations

Training data from multiple sources



Training data from multiple sources



Person sleeping at desk icon made by Freepik
from www.flaticon.com

How much can be learned even if some data is corrupted or manipulated?

Overview

Reminder: (Statistical) Learning Theory

Robust Learning From Untrusted Sources

Robust Fair Learning

Slides available at: <http://cvm1.ist.ac.at>

Reminder: Supervised Learning

Setting:

- ▶ **Inputs:** $x \in \mathcal{X}$, e.g. strings, images, vectors, ...
- ▶ **Outputs:** $y \in \mathcal{Y}$. For simplicity, we use $\mathcal{Y} = \{\pm 1\}$ (binary classification)
- ▶ **Probability distribution:** $p(x, y)$ over $\mathcal{X} \times \mathcal{Y}$, unknown to the learner
- ▶ **Loss function:** $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. For simplicity, we use 0/1-loss: $\ell(y, \bar{y}) = \mathbb{I}[y \neq \bar{y}]$

Abstract Goal:

- ▶ find a **prediction function**, $f : \mathcal{X} \rightarrow \mathcal{Y}$, such that the expected loss

$$\text{er}(h) = \mathbb{E}_{(x, y) \sim p}[\ell(y, f(x))] = \Pr_{(x, y) \sim p}\{f(x) \neq y\}$$

on *future data* is small.

Learning from data:

- ▶ training data: $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \stackrel{i.i.d.}{\sim} p$
- ▶ hypothesis class: $\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$
- ▶ learning algorithm $\mathcal{L} : \mathbb{P}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{H}$, $\mathbb{P}(\cdot) = \text{power set}$
 - ▶ input: a training set, $S \subset \mathcal{X} \times \mathcal{Y}$,
 - ▶ output: a trained model $\mathcal{L}(S) \in \mathcal{H}$ (= prediction function).

Learning from data:

- ▶ training data: $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \stackrel{i.i.d.}{\sim} p$
- ▶ hypothesis class: $\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$
- ▶ learning algorithm $\mathcal{L} : \mathbb{P}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{H}$, $\mathbb{P}(\cdot) =$ power set
 - ▶ input: a training set, $S \subset \mathcal{X} \times \mathcal{Y}$,
 - ▶ output: a trained model $\mathcal{L}(S) \in \mathcal{H}$ (= prediction function).

Central question in Statistical Learning Theory:

Is there a universal learning algorithm, such that: $\text{er}(\mathcal{L}(S)) \stackrel{|S| \rightarrow \infty}{\rightarrow} \min_{h \in \mathcal{H}} \text{er}(h)$?

Learning from data:

- ▶ training data: $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \stackrel{i.i.d.}{\sim} p$
- ▶ hypothesis class: $\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$
- ▶ learning algorithm $\mathcal{L} : \mathbb{P}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{H}$, $\mathbb{P}(\cdot) =$ power set
 - ▶ input: a training set, $S \subset \mathcal{X} \times \mathcal{Y}$,
 - ▶ output: a trained model $\mathcal{L}(S) \in \mathcal{H}$ (= prediction function).

Central question in Statistical Learning Theory:

Is there a universal learning algorithm, such that: $\text{er}(\mathcal{L}(S)) \xrightarrow{|S| \rightarrow \infty} \min_{h \in \mathcal{H}} \text{er}(h)$?

Classic result: [Vapnik&Chervonenkis, 1971], [Blumer, Ehrenfeucht, Hassler, Warmuth, 1989]

If and only if $\text{VC}(\mathcal{H}) < \infty$, empirical risk minimization (ERM) does the job:

$$\mathcal{L}(S) \leftarrow \underset{h \in \mathcal{H}}{\text{argmin}} \text{er}_S(h) \quad \text{for } \text{er}_S(h) := \frac{1}{|S|} \sum_{(x,y) \in S} \mathbb{I}[f(x) \neq y].$$

[V. N. Vapnik, A. Ya. Chervonenkis. "Theory of uniform convergence of frequencies of appearance of attributes to their probabilities and problems of defining optimal solution by empiric data". Theory of Probability and its Applications, 1971]

[A. Blumer, A. Ehrenfeucht, D. Haussler, M. K. Warmuth. "Learnability and the Vapnik-Chervonenkis Dimension". Journal of the ACM, 1989]

Learning from unreliable/malicious data:

- ▶ training set: $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$
- ▶ but: data has issues: some data points might not really be samples from p

Learning from unreliable/malicious data:

- ▶ training set: $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$
- ▶ but: data has issues: some data points might not really be samples from p
- ▶ formally: malicious adversary \mathcal{A} [Valiant, 1985]
 - ▶ \mathcal{A} can manipulate a fraction α of the dataset
 - ▶ input: dataset S
 - ▶ output: dataset $S' = \mathcal{A}(S)$ with $\lceil (1 - \alpha)m \rceil$ points are unchanged and $\lfloor \alpha m \rfloor$ are arbitrary
 - ▶ \mathcal{A} can depend on the learning algorithms, etc.

Question: Is ERM still be a universally good learning strategy?

Learning from unreliable/malicious data:

- ▶ training set: $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$
- ▶ but: data has issues: some data points might not really be samples from p
- ▶ formally: malicious adversary \mathcal{A} [Valiant, 1985]
 - ▶ \mathcal{A} can manipulate a fraction α of the dataset
 - ▶ input: dataset S
 - ▶ output: dataset $S' = \mathcal{A}(S)$ with $\lceil(1 - \alpha)m\rceil$ points are unchanged and $\lfloor\alpha m\rfloor$ are arbitrary
 - ▶ \mathcal{A} can depend on the learning algorithms, etc.

Question: Is ERM still be a universally good learning strategy?

Classic Result: no! [Kearns and Li, 1993]

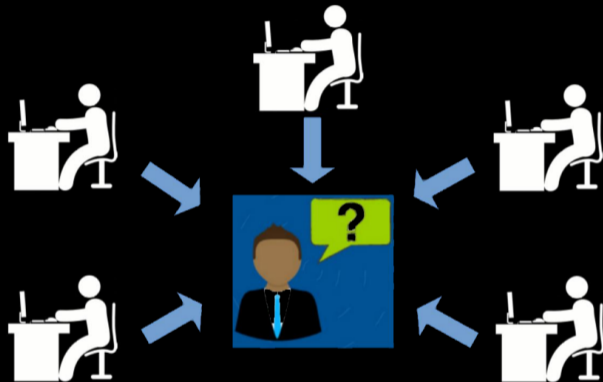
No learning algorithm can guarantee an error less than $\frac{\alpha}{1-\alpha}$ on future data!

[L. G. Valiant. "Learning disjunctions of conjunctions". IJCAI 1985]

[M. Kearns, M. Li. "Learning in the presence of malicious errors". SIAM Journal on Computing, 1993]

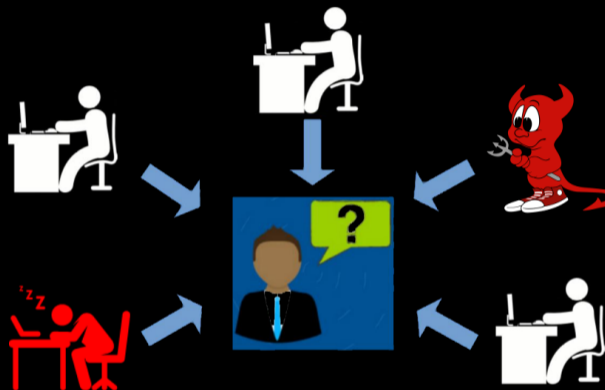
Learning from Multiple Sources

Training data from multiple sources



If all sources are i.i.d. samples from the correct data distribution
→ naive strategy "merge all datasets and train a classifier" works perfectly

Training data from multiple sources



Person sleeping at desk icon made by Freepik
from www.flaticon.com

If some sources are not reliable, naive strategy can fail miserably!

Robust Learning from Unreliable or Malicious Sources



Nikola
Konstantinov



Elias
Frantar



Dan
Alistarh

Disclaimer: "These results have been modified from their original form. They have been edited to fit the screen and the allotted time slot."

[N. Konstantinov, E. Frantar, D. Alistarh, CHL. "On the Sample Complexity of Adversarial Multi-Source PAC Learning", ICML 2020]

[N. Konstantinov, CHL. "Robust Learning from Untrusted Sources", ICML 2019]

Learning from Multiple Sources

- ▶ multiple training sets: S_1, S_2, \dots, S_N
 - ▶ each $S_i = \{(x_1^i, y_1^i), \dots, (x_m^i, y_m^i)\} \stackrel{i.i.d.}{\sim} p$
- ▶ multi-source learning algorithm: $\mathcal{L} : (\mathcal{X} \times \mathcal{Y})^{N \times m} \rightarrow \mathcal{H}$
 - ▶ input: training sets, S_1, S_2, \dots, S_N
 - ▶ output: one hypothesis $\mathcal{L}(S_1, \dots, S_N) \in \mathcal{H}$ (= a trained model).

Learning from Multiple Unreliable/Malicious Sources

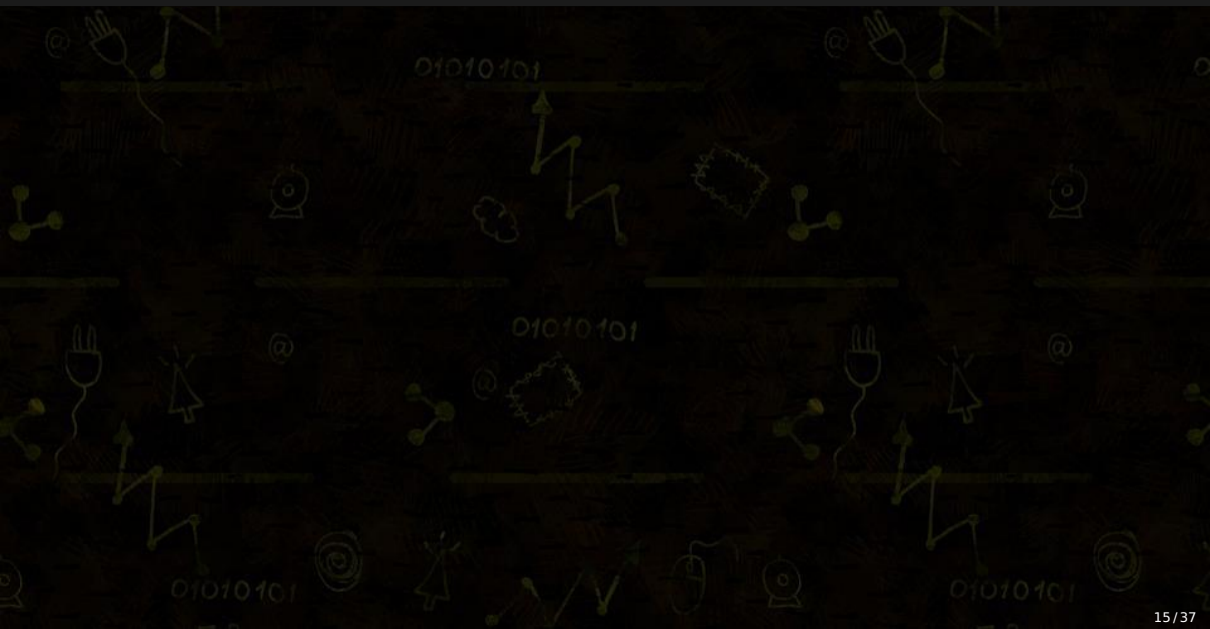
- ▶ multiple training sets: S_1, S_2, \dots, S_N
 - ▶ each $S_i = \{(x_1^i, y_1^i), \dots, (x_m^i, y_m^i)\} \stackrel{i.i.d.}{\sim} p$
- ▶ multi-source learning algorithm: $\mathcal{L} : (\mathcal{X} \times \mathcal{Y})^{N \times m} \rightarrow \mathcal{H}$
 - ▶ input: training sets, $S'_1, S'_2, \dots, S'_N = \mathfrak{A}(S_1, \dots, S_N)$
 - ▶ output: one hypothesis $\mathcal{L}(S'_1, S'_2, \dots, S'_N) \in \mathcal{H}$ (= a trained model).
- ▶ adversary \mathfrak{A}
 - ▶ input: data sets S_1, \dots, S_N
 - ▶ output: data sets S'_1, \dots, S'_N , of which $\lceil (1 - \alpha)N \rceil$ are identical to before and $\lfloor \alpha N \rfloor$ are arbitrary
 - ▶ the adversary might know the training algorithm

Learning from Multiple Unreliable/Malicious Sources

- ▶ multiple training sets: S_1, S_2, \dots, S_N
 - ▶ each $S_i = \{(x_1^i, y_1^i), \dots, (x_m^i, y_m^i)\} \stackrel{i.i.d.}{\sim} p$
- ▶ multi-source learning algorithm: $\mathcal{L} : (\mathcal{X} \times \mathcal{Y})^{N \times m} \rightarrow \mathcal{H}$
 - ▶ input: training sets, $S'_1, S'_2, \dots, S'_N = \mathfrak{A}(S_1, \dots, S_N)$
 - ▶ output: one hypothesis $\mathcal{L}(S'_1, S'_2, \dots, S'_N) \in \mathcal{H}$ (= a trained model).
- ▶ adversary \mathfrak{A}
 - ▶ input: data sets S_1, \dots, S_N
 - ▶ output: data sets S'_1, \dots, S'_N , of which $\lceil (1 - \alpha)N \rceil$ are identical to before and $\lfloor \alpha N \rfloor$ are arbitrary
 - ▶ the adversary might know the training algorithm

Is there a universal learning algorithm, i.e. $\text{er}(\mathcal{L}(S'_1, \dots, S'_N)) \xrightarrow{m \rightarrow \infty} \min_{h \in \mathcal{H}} \text{er}(h)$?

Related Work



Robust learning from a single dataset

- ▶ no universal algorithm: minimum guaranteeable error is $\frac{\alpha}{1-\alpha}$ [Kearns and Li, 1993]
- ▶ identical to our situation when each dataset consists of a single point, $m = 1$
 - only $N \rightarrow \infty$ will probably not suffice to learn arbitrarily well

Robust learning from a single dataset

- ▶ no universal algorithm: minimum guaranteeable error is $\frac{\alpha}{1-\alpha}$ [Kearns and Li, 1993]
- ▶ identical to our situation when each dataset consists of a single point, $m = 1$
 - only $N \rightarrow \infty$ will probably not suffice to learn arbitrarily well

Collaborative learning (multiple parties together learn *individual models*)

- ▶ universal learning algorithm exists [Blum et al., 2017], [Qiao, 2018]

Robust learning from a single dataset

- ▶ no universal algorithm: minimum guaranteeable error is $\frac{\alpha}{1-\alpha}$ [Kearns and Li, 1993]
- ▶ identical to our situation when each dataset consists of a single point, $m = 1$
→ only $N \rightarrow \infty$ will probably not suffice to learn arbitrarily well

Collaborative learning (multiple parties together learn *individual models*)

- ▶ universal learning algorithm exists [Blum et al., 2017], [Qiao, 2018]

Density estimation from untrusted batches

- ▶ possible, but not applicable to supervised learning [Qiao and Valiant, 2018], [Jain and Orlicsky, 2020]

Robust learning from a single dataset

- ▶ no universal algorithm: minimum guaranteeable error is $\frac{\alpha}{1-\alpha}$ [Kearns and Li, 1993]
- ▶ identical to our situation when each dataset consists of a single point, $m = 1$
→ only $N \rightarrow \infty$ will probably not suffice to learn arbitrarily well

Collaborative learning (multiple parties together learn *individual models*)

- ▶ universal learning algorithm exists [Blum et al., 2017], [Qiao, 2018]

Density estimation from untrusted batches

- ▶ possible, but not applicable to supervised learning [Qiao and Valiant, 2018], [Jain and Orlicsky, 2020]

Byzantine-robust distributed optimization

- ▶ specific solutions for gradient-based optimization [Yin et al., 2018], [Alistarh et al., 2018]
- ▶ results focus on convergence analysis

Theorem [Konstantinov et al., 2020]

There exists a learning algorithm, \mathcal{L} , such that with high probability:

$$\text{er}(\mathcal{L}(S'_1, \dots, S'_N)) \leq \min_{h \in \mathcal{H}} \text{er}(h) + \underbrace{\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{(1-\alpha)Nm}} + \alpha \frac{1}{\sqrt{m}}\right)}_{\rightarrow 0 \text{ for } m = |S| \rightarrow \infty},$$

with $S'_1, \dots, S'_N = \mathcal{A}(S_1, \dots, S_N)$ for any adversary \mathcal{A} with $\alpha < \frac{1}{2}$.

($\tilde{\mathcal{O}}$ -notation hides constant and logarithmic factors)

Big Picture

01010101

Question: why is learning easier from multiple sources than from a single one?

Answer: it's not. But the task for the adversary is harder!

- ▶ single source: no restrictions how to manipulate the data
- ▶ multi-source: manipulation must adhere to the source structure

Algorithm idea: exploit law of large numbers

1. majority of datasets are unperturbed
2. for $m \rightarrow \infty$ these start to look more and more similar
3. we can identify (at least) the unperturbed datasets
4. we perform ERM on the union of only those

01010101

01010101

Robust multi-source learning algorithm:

- ▶ **Input:** datasets S'_1, \dots, S'_N
- ▶ **Input:** suitable distance measure d between datasets
- ▶ **Input:** suitable threshold value θ
- ▶ Step 1) identify which sources to trust
 - ▶ compute all pairwise distance d_{ij} between datasets S'_1, \dots, S'_N
 - ▶ for any i : if $d_{ij} < \theta$ for at least $\lfloor \frac{N}{2} \rfloor$ values of $j \neq i$, then $T \leftarrow T \cup \{i\}$
- ▶ Step 2) merge data from all sources S'_i with $i \in T$ into a new dataset \tilde{S}
- ▶ Step 3) minimize training error on \tilde{S}

Open choices:

- ▶ distance measure d (discussed later), threshold θ (see paper)



All datasets clean



All datasets clean



All datasets clean



All datasets clean



All datasets clean



All datasets clean



All datasets clean



All datasets clean



All datasets clean → all datasets included → same as (optimal) naive algorithm



Some datasets manipulated



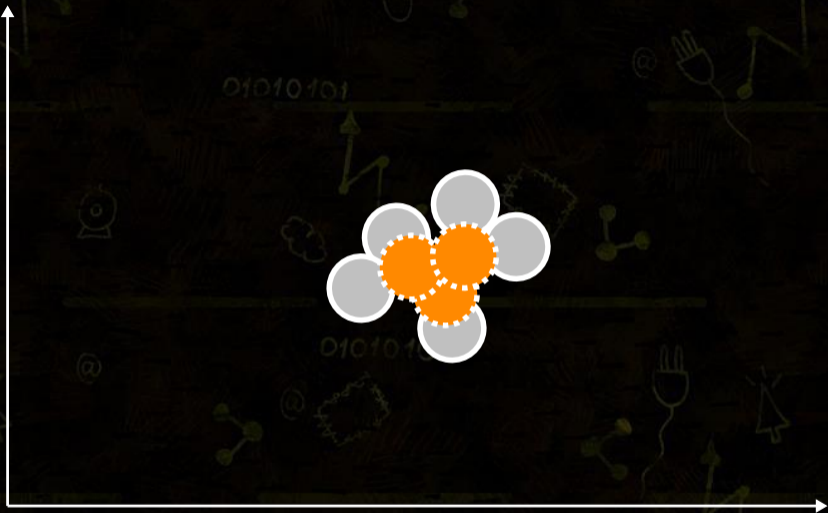
Some datasets manipulated → manipulated datasets excluded



Consistent manipulations



Consistent manipulations → manipulated datasets excluded



Some datasets manipulated to look like originals



Some datasets manipulated to look like originals → all datasets included.

What properties does the distance measure d need?

01010101

01010101

01010101

01010101

What properties does the distance measure d need?

1) 'clean' datasets should get grouped together:

$$S, \hat{S} \sim p \Rightarrow d(S, \hat{S}) \xrightarrow{m \rightarrow \infty} 0$$

What properties does the distance measure d need?

1) 'clean' datasets should get grouped together:

$$S, \hat{S} \sim p \Rightarrow d(S, \hat{S}) \xrightarrow{m \rightarrow \infty} 0$$

2) if manipulated datasets are grouped with the clean ones, they should not hurt the learning step

$$d(S, \hat{S}) \text{ is small} \Rightarrow \mathcal{L}(\hat{S}) \approx \mathcal{L}(S)$$

What properties does the distance measure d need?

1) 'clean' datasets should get grouped together:

$$S, \hat{S} \sim p \Rightarrow d(S, \hat{S}) \xrightarrow{m \rightarrow \infty} 0$$

2) if manipulated datasets are grouped with the clean ones, they should not hurt the learning step

$$d(S, \hat{S}) \text{ is small} \Rightarrow \mathcal{L}(\hat{S}) \approx \mathcal{L}(S)$$

Observation:

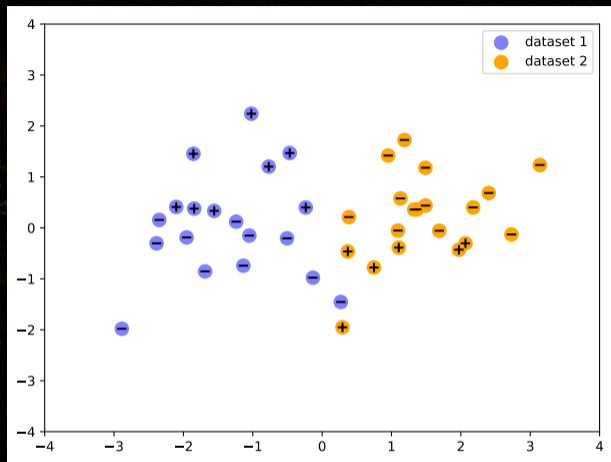
- ▶ many candidate distances do not fulfill both conditions simultaneously:
 - ▶ geometric: average Euclidean distance, Chamfer distance, Hausdorff distance, ...
 - ▶ probabilistic: Wasserstein distance, total variation, KL-divergence, ...
- ▶ **discrepancy distance** does fulfill the conditions!

Discrepancy Distance [Mansour et al. 2009]

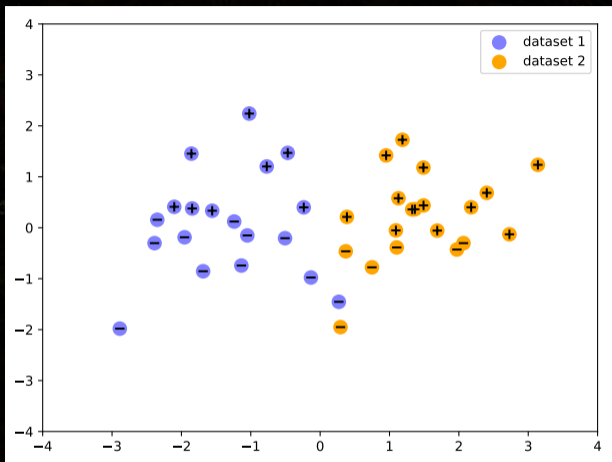
For a set of classifiers \mathcal{H} and datasets S, \hat{S} , define

$$\text{disc}(S, \hat{S}) = \max_{h \in \mathcal{H}} |\text{er}_S(h) - \text{er}_{\hat{S}}(h)|.$$

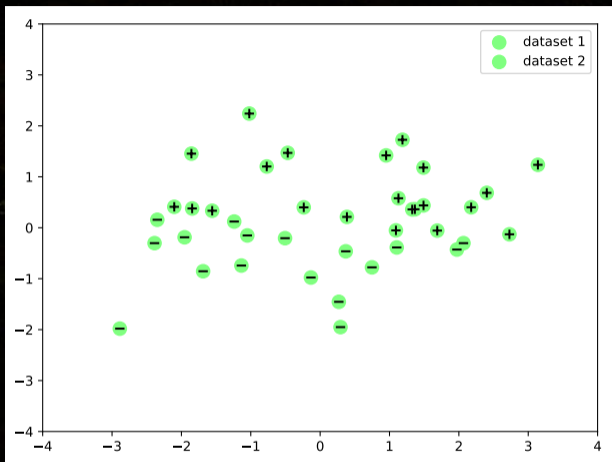
- ▶ maximal amount any classifier, $h \in \mathcal{H}$, can disagree between S, \hat{S}
- ▶ discrepancy can be estimated by training a classifier itself:
 - ▶ $S^\pm \leftarrow S$ with all ± 1 labels flipped to their opposites
 - ▶ $\tilde{S} \leftarrow S^\pm \cup \hat{S}$
 - ▶ $\text{disc}(S, \hat{S}) \leftarrow 1 - 2 \min_{h \in \mathcal{H}} \text{er}_{\tilde{S}}(h)$ (minimal training error of any $h \in \mathcal{H}$ on \tilde{S})



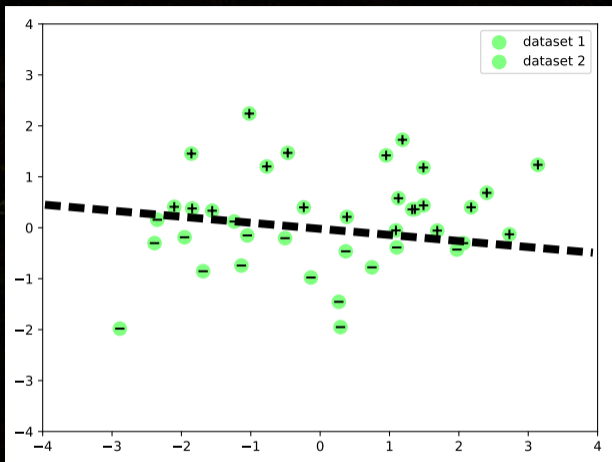
Two datasets, S, \hat{S}



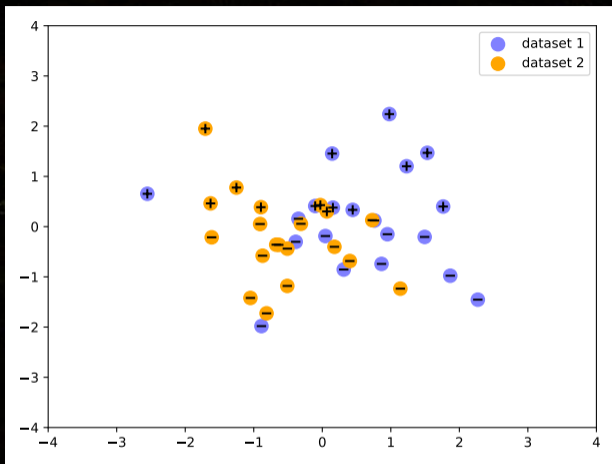
Flip signs of S



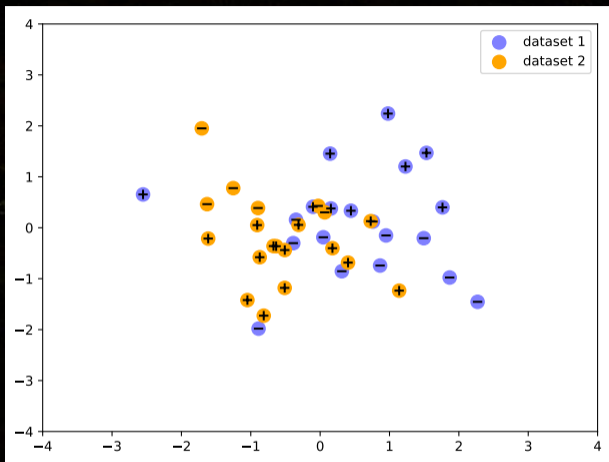
Merge both datasets



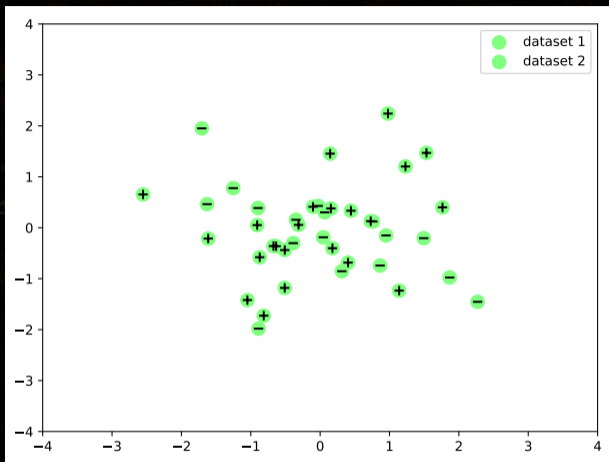
Classifier with small training error \rightarrow large discrepancy



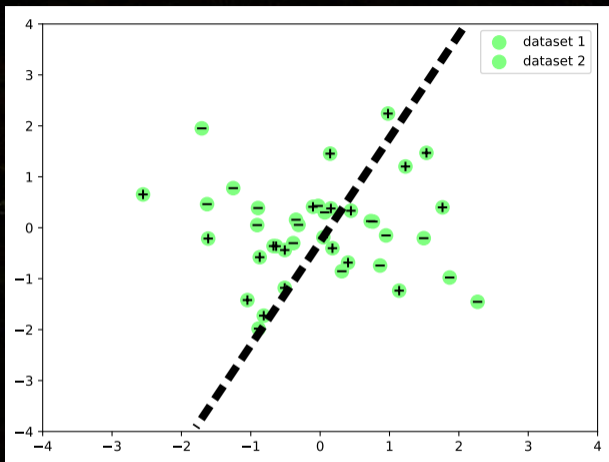
Two datasets, S, \hat{S}



Flip signs of S



Merge both datasets



No classifier with small training error \rightarrow small discrepancy

Observation: discrepancy distance has both property we need

- 1) Datasets from the same distribution (eventually) get grouped together
 - ▶ for $\mathbf{VC}(\mathcal{H}) < \infty$, if S and \hat{S} are sampled from the same distribution, then

$$\text{disc}(S, \hat{S}) \rightarrow 0 \quad \text{for} \quad |S|, |\hat{S}| \rightarrow \infty$$

- 2) Datasets that are grouped together cannot hurt the learning much

Consider:

- ▶ training set $S_{\text{trn}} \stackrel{i.i.d.}{\sim} p$
- ▶ arbitrary set \hat{S} , potentially manipulated but with $\text{disc}(S_{\text{trn}}, \hat{S}) \leq \theta$
- ▶ test set $S_{\text{tst}} \stackrel{i.i.d.}{\sim} p$

Then, for every $h \in \mathcal{H}$:

$$\text{er}_{S_{\text{tst}}}(h) \leq \text{er}_{\hat{S}}(h) + \underbrace{\text{disc}(S_{\text{trn}}, \hat{S})}_{\leq \theta} + \underbrace{\text{disc}(S_{\text{trn}}, S_{\text{tst}})}_{\text{small by prop. 1}}$$

Robust Fair Learning

Fairness-Aware Learning from Unreliable or Malicious Data



Nikola
Konstantinov



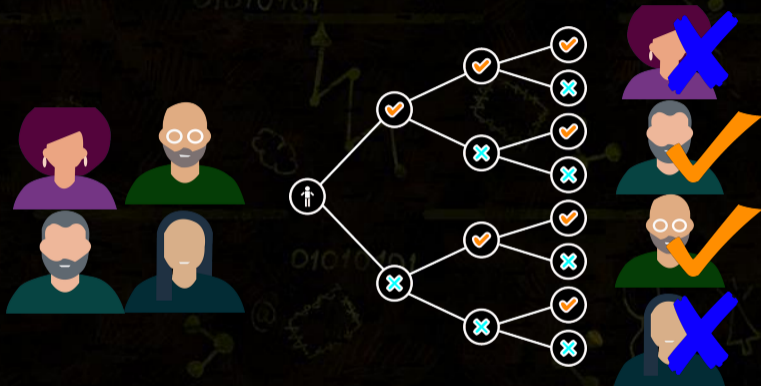
Jen
Iofinova

Disclaimer: "These results have been modified from their original form. They have been edited to fit the screen and the allotted time slot."

[N. Konstantinov, CHL. "*Fairness-Aware PAC Learning from Corrupted Data*", <https://arxiv.org/abs/2102.06004>]

[E. Iofinova*, N. Konstantinov*, CHL. "*Robust Learning from Untrusted Sources*", <https://arxiv.org/abs/2106.11732>]

Algorithmic Fairness



How to ensure that a classifier does not discriminate against certain groups?

Setting:

- ▶ **Inputs:** $x \in \mathcal{X}$, e.g. strings, images, vectors, ...
- ▶ **Protected attribute:** $a \in \mathcal{A}$, e.g. gender, age, race, ...
- ▶ **Outputs:** $y \in \mathcal{Y}$ (for simplicity: $\mathcal{Y} = \{0, 1\}$)
- ▶ **Probability distribution:** $p(x, a, y)$ over $\mathcal{X} \times \mathcal{A} \times \mathcal{Y}$
- ▶ **Loss function:** $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ (for simplicity: 0/1-loss)

Abstract Goal:

- ▶ find a **prediction function**, $f : \mathcal{X} \rightarrow \mathcal{Y}$ low expected loss

$$\text{er}(h) = \mathbb{E}_{(x,y) \sim p}(\mathbb{1}[f(x) \neq y]) = \Pr_{(x,y) \sim p}\{f(x) \neq y\}$$

that in addition **fulfills some condition of (group) fairness.**

Group Fairness:

- ▶ **demographic parity:** "all groups have the same success rate"

$$\forall a, b \in \mathcal{A} \quad p(f(X) = 1 | A = a) = p(f(X) = 1 | A = b)$$

- ▶ **equality of opportunity:** "all groups have the true positive rate"

$$\forall a, b \in \mathcal{A} \quad p(f(X) = 1 | A = a, Y = 1) = p(f(X) = 1 | A = b, Y = 1)$$

and many others [Barocas et al., 2019]

Several fairness-aware learning methods exist to enforce these criteria.

Fair Learning from unreliable/malicious data:

- ▶ original training set: $S = \{(x_1, a_1, y_1), \dots, (x_m, a_m, y_m)\}$
- ▶ adversary \mathcal{A} can manipulate a fraction α of the dataset
- ▶ actual training set: $\mathcal{A}(S)$

Question: Can a fairness-aware learner overcome the manipulation?

Fair Learning from unreliable/malicious data:

- ▶ original training set: $S = \{(x_1, a_1, y_1), \dots, (x_m, a_m, y_m)\}$
- ▶ adversary \mathcal{A} can manipulate a fraction α of the dataset
- ▶ actual training set: $\mathcal{A}(S)$

Question: Can a fairness-aware learner overcome the manipulation?

Theorem [Konstantinov and Lampert, 2021]

There is even for finite-sized hypothesis classes, \mathcal{H} , for which:

- ▶ No learning algorithm can guarantee optimal fairness.
- ▶ This effect is independent of whether accuracy is also affected or not.
- ▶ The smaller the minority group, the stronger the bias.

Fairness-Aware Learning from Multiple Unreliable Sources

- ▶ multiple training sets: $S_1, S_2, \dots, S_N \subset \mathcal{X} \times \mathcal{A} \times \mathcal{Y}$
- ▶ adversary \mathfrak{A} can manipulate $K = \lfloor \alpha N \rfloor$ of the datasets for $\alpha < \frac{1}{2}$
- ▶ actual training sets: $\mathfrak{A}(S_1, \dots, S_N)$

Is there a fairness-aware learning algorithm that overcomes such manipulations?

Fairness-Aware Learning from Multiple Unreliable Sources

- ▶ multiple training sets: $S_1, S_2, \dots, S_N \subset \mathcal{X} \times \mathcal{A} \times \mathcal{Y}$
- ▶ adversary \mathfrak{A} can manipulate $K = \lfloor \alpha N \rfloor$ of the datasets for $\alpha < \frac{1}{2}$
- ▶ actual training sets: $\mathfrak{A}(S_1, \dots, S_N)$

Is there a fairness-aware learning algorithm that overcomes such manipulations?

Theorem [Iofinova et al., 2021]

There exists a learning algorithm, \mathcal{L} , such that for $h^* = \mathcal{L}(\mathfrak{A}(S_1, \dots, S_N))$ with high probability

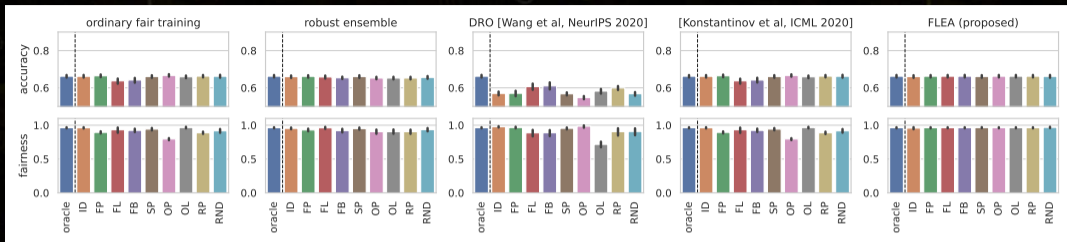
$$\text{er}(h^*) \leq \min_{h \in \mathcal{H}} \text{er}(h) + \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{m}}\right), \quad \Gamma(h^*) \leq \min_{h \in \mathcal{H}} \Gamma(h) + \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{m}}\right)$$

where Γ is a quantitative measure of *demographic parity* fairness.

FLEA (Fair LEarning against Adversaries):

- ▶ **Input:** datasets S'_1, \dots, S'_N
- ▶ **Input:** $\beta \leq \frac{1}{2}$ upper bound on fraction of malignant sources
- ▶ **Define:** distance measure $d(S, \hat{S}) = \text{disc}(S, \hat{S}) + \text{disp}(S, \hat{S}) + \text{disb}(S, \hat{S})$
 - ▶ $\text{disc}(S, \hat{S})$: discrepancy as before
 - ▶ $\text{disp}(S, \hat{S})$: maximal fairness difference of any classifier between S and \hat{S}
 - ▶ $\text{disb}(S, \hat{S})$: difference in protected group proportions
- ▶ Step 1) identify which sources to trust
 - ▶ compute all pairwise distance d_{ij} between datasets S'_1, \dots, S'_N
 - ▶ for any $i = 1, \dots, N$: $q_i \leftarrow \beta$ -quantile(d_{i1}, \dots, d_{iN})
 - ▶ $T \leftarrow \{i : q_i \leq \beta\text{-quantile}(q_1, \dots, q_N)\}$
- ▶ Step 2) merge data from all sources S'_i with $i \in T$ into a new dataset \tilde{S}
- ▶ Step 3) train fairness-aware learning algorithm on \tilde{S}

Experimental Results



- ▶ bars are different data manipulations, designed to hurt accuracy or fairness
- ▶ simply training on all data often suboptimal
- ▶ other baselines often fail to overcome problems
- ▶ FLEA reliably recovers fairness and accuracy

method	COMPAS	
	accuracy	fairness
naïve	63.5 \pm 2.1	78.9 \pm 2.3
robust ensemble	65.0 \pm 1.1	88.4 \pm 2.9
DRO (Wang et al., 2020)	54.5 \pm 1.2	70.9 \pm 5.7
(Konstantinov et al., 2020)	63.5 \pm 2.1	78.9 \pm 2.3
FLEA (proposed)	65.9 \pm 1.1	95.3 \pm 2.3
oracle	66.2 \pm 1.1	96.2 \pm 1.3

More results and ablation studies in the paper.

Summary

Bad news:

- ▶ Learning is not robust to bad data.
- ▶ This can affect accuracy as well as fairness.

Good news:

- ▶ Modern data set are often not monolithic but collected from multiple sources.
- ▶ Multi-source learning **can** be made robust to bad data sources.
- ▶ This holds for accuracy as well as fairness.

Thank you!

Thanks to:



Nikola Konstantinov



Jen Iofinova



Elias Frantar



Dan Alistarh

Funding sources:



References

- D. Alistarh, Z. Allen-Zhu, and J. Li. Byzantine stochastic gradient descent. In *NeurIPS*, 2018.
- S. Barocas, M. Hardt, and A. Narayanan. *Fairness and Machine Learning*. fairmlbook.org, 2019. <http://www.fairmlbook.org>.
- A. Blum, N. Haghtalab, A. D. Procaccia, and M. Qiao. Collaborative pac learning. In *NIPS*. 2017.
- E. Iofinova, N. Konstantinov, and C. H. Lampert. FLEA: Provably robust fair multisource learning. *arXiv: 2106.11732 [cs.LG]*, 2021.
- A. Jain and A. Orlitsky. Optimal robust learning of discrete distributions from batches. In *ICML*, 2020.
- M. Kearns and M. Li. Learning in the presence of malicious errors. In *SIAM Journal on Computing*, 1993.
- N. Konstantinov and C. H. Lampert. Robust learning from untrusted sources. In *ICML*, 2019.
- N. Konstantinov and C. H. Lampert. Fairness-aware PAC learning from corrupted data. *arXiv: 2102.06004 [cs.LG]*, 2021.
- N. Konstantinov, E. Frantar, D. Alistarh, and C. H. Lampert. On the sample complexity of adversarial multi-source PAC learning. In *ICML*, 2020.
- M. Qiao. Do outliers ruin collaboration? In *ICML*, 2018.
- M. Qiao and G. Valiant. Learning discrete distributions from untrusted batches. In *ITCS*, 2018.
- L. G. Valiant. Learning disjunctons of conjunctions. In *IJCAI*, 1985.
- D. Yin, Y. Chen, K. Ramchandran, and P. Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *ICML*, 2018.