# Scalable and Provably Robust Algorithms for Machine Learning



### Yu Cheng

University of Illinois at Chicago

# Motivation

#### The Washington Post Democracy Dies in Darkness

#### WorldViews

# Syrian hackers claim AP hack that tipped stock market by \$136 billion.

By **Max Fisher** April 23, 2013





Can we design algorithms that are robust when a small fraction of the data is adversarially corrupted?

# Motivation



TECHNOLOGY

#### Researchers Develop Glasses Frames That Fool Facial Recognition Technology

November 5, 2016 · 8:27 AM ET Heard on Weekend Edition Saturday















Yu Cheng



4 Aug 2017 | 18:00 GMT

#### Slight Street Sign Modifications Can Completely Fool Machine Learning Algorithms

Minor changes to street sign graphics can fool machine learning algorithms into thinking the signs say something completely different







[Eykholt+ '18]

Can we make non-convex optimization robust?

# Motivation

#### Vex

#### recode

### Artificial intelligence will help determine if you get your next job

Al is being used to attract applicants and to predict a candidate's fit for a position. But is it up to the task? By Rebecca Heilweil | Dec 12, 2019, 8:00am EST Technology | Future Finance

# Sweet-Talking CEOs Are Starting to Outsmart the Robot Analysts

**Bloomberg** 

Study finds companies alter words to cater to listening algos
 Emphasis on positivity as negative phrases get ditched

By <u>Gregor Stuart Hunter</u> +Follow October 20, 2020, 8:03 AM EDT

# Can we design ML systems that work well when the data is provided by self-interested agents?

Ξ

Subscribe

Challenges

#### What are challenges in developing robust ML algorithms?

Challenges

#### What are challenges in developing robust ML algorithms?



Yu Cheng

Challenges

#### What are challenges in developing robust ML algorithms?

• Adversarial attacks.





# Challenges

What are challenges in developing robust ML algorithms?

- Adversarial attacks.
- Model misspecification/Strong assumptions.





# Challenges

What are challenges in developing robust ML algorithms?

- Adversarial attacks.
- Model misspecification/Strong assumptions.
- Strategic behaviors.



# Challenges

What are challenges in developing robust ML algorithms?

- Adversarial attacks.
- Model misspecification/Strong assumptions.
- Strategic behaviors.

Design fast and provably robust algorithms for ML.

- Focus on basic problems that are well-understood w/o corruption.
- Bring together ideas from ML, optimization, and game theory.



Mechanism Design and Equilibrium Computation [CGMW WINE'18 Best Paper] [CDS SODA'17] [CCT ITCS'17]

> Graph Sparsification [CCPS ICALP'21] [CCLPT COLT'15]

Information Structure Design [BCKS EC'16] [CCDEHT FOCS'15]

Fairness and Social Choice [KJWCM AISTATS'21] **[C**]MW EC'19] [ZCC AAAI'19] [CDK AAAI'18] [CDK EC'17]

Planning in MDPs with Agents [ZCC AAAI'22] [ZCC AAAI'21]

# Outline

Part I: Introduction

• Motivation/Challenges

#### Part II: Robust Algorithms for ML

- High-Dimensional Statistics
- Non-Convex Optimization
- Learning with Strategic Agents

#### Part III: Future Directions



## Mean Estimation

- Input: *n* samples  $\{X_1, \dots, X_n\}$  drawn from  $\mathcal{N}(\mu^*, I)$  on  $\mathbb{R}^d$ .
- Goal: Learn  $\mu^*$ .

Empirical mean 
$$\widehat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i$$
 works:  
 $\|\widehat{\mu} - \mu^*\|_2 \le \epsilon$  when  $n = \Omega(d/\epsilon^2)$ .

## **Robust** Mean Estimation



# $\epsilon$ -Corruption Model



### Goal: given an $\epsilon$ -corrupted set of n samples, learn $\mu^*$ .

# Classical Estimators and Why They Fail

- Empirical mean
- Empirical median
  - Coordinate-wise median: can be  $\Theta(\epsilon \sqrt{d})$  far from the true mean.
  - Tukey median: O(ε) error but is
    NP-Hard to compute.



# **Robust Mean Estimation**

Algorithm	Error Guarantee	Runtime					
Coordinate-wise median	$O(\epsilon \sqrt{d})$	0(nd)					
Tukey median	$O(\epsilon)$	NP-Hard					

# **Robust Mean Estimation**

Algorithm	Error Guarantee	Runtime						
Coordinate-wise median	$O(\epsilon \sqrt{d})$	0(nd)						
Tukey median	$O(\epsilon)$	NP-Hard						
[Lai+ '16]	$O(\epsilon \sqrt{\log d})$	D 1 1						
[Diakonikolas+ '16]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	Polynomial						

There has been a flurry of research that obtained polynomial-time robust algorithms for a wide range of high-dimensional learning tasks.

# High-Dimensional Robust Statistics

When a small fraction of the input is corrupted:

[Tukey '60, Huber '64]: provably robust statistical estimators.



# High-Dimensional Robust Statistics

When a small fraction of the input is corrupted:

[Tukey '60, Huber '64]: provably robust statistical estimators.

[Diakonikolas+ '16, Lai+ '16]: provably robust statistical estimators that can be computed in **polynomial-time**.

Polynomial-time ≠ Scalable!

These algorithms are often much slower than the fastest non-robust ones.

# High-Dimensional Robust Statistics

When a small fraction of the input is corrupted:

[Tukey '60, Huber '64]: provably robust statistical estimators.

[Diakonikolas+ '16, Lai+ '16]: provably robust statistical estimators that can be computed in **polynomial-time**.

[C Diakonikolas Ge '19]: provably robust statistical estimators that are **as efficient as** their non-robust counterparts.

# **Robust Mean Estimation**

Algorithm	Error Guarantee	Runtime					
Coordinate-wise median	$O(\epsilon \sqrt{d})$	<i>O</i> ( <i>nd</i> )					
Tukey median	$O(\epsilon)$	NP-Hard					
[Lai+ '16]	$O(\epsilon \sqrt{\log d})$	Dolumomial					
[Diakonikolas+ '16]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	Polynomiai					

# **Robust Mean Estimation**

Algorithm	Error Guarantee	Runtime						
Coordinate-wise median	$O(\epsilon \sqrt{d})$	0(nd)						
Tukey median	$O(\epsilon)$	NP-Hard						
[Lai+ '16]	$O(\epsilon \sqrt{\log d})$	$O(md^2)$						
[Diakonikolas+ '16]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	<u>sz(nu</u> )						

### Robust mean estimation in nearly-linear time?

# Robust Mean Estimation in Nearly Linear Time [C Diakonikolas Ge '19]

Algorithm	Error Guarantee	Runtime						
Coordinate-wise median	$O(\epsilon \sqrt{d})$	<i>O</i> ( <i>nd</i> )						
Tukey median	$O(\epsilon)$	NP-Hard						
[Lai+ '16]	$O(\epsilon \sqrt{\log d})$	$O(md^2)$						
[Diakonikolas+ '16]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	$SZ(na^{-})$						
[ <b>C</b> Diakonikolas Ge '19]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	$ ilde{O}(nd)/ ext{poly}(\epsilon)$						

For constant  $\epsilon$ , our algorithm has the best-possible error guarantee, sample complexity, and running time (up to log factors).

# Robust Mean Estimation in Nearly Linear Time [C Diakonikolas Ge '19]



With corruption, the top eigenvector of M can be an arbitrary vector y'!

#### Robust Mean Estimation in Nearly Linear Time [C Diakonikolas Ge '19]



#### Robust Mean Estimation in Nearly Linear Time [C Diakonikolas Ge '19]



Yu Cheng

# Robust Mean Estimation in Nearly Linear Time [C Diakonikolas Ge '19]

- Start with coordinate-wise median that is  $O(\sqrt{d})$  far from  $\mu^*$ .
- The distance to  $\mu^*$  decreases by half in each iteration.
  - Each iteration takes nearly-linear time (mostly solving a packing SDP).



# Beyond Mean Estimation

provably robust statistical estimators that are **as efficient as** their non-robust counterparts.

[**C**DG'19]:

[DL'19][DHL'19]:

[**C**DGW'19][LY'20]:

[C+'20]:

[C+'20][D+'20]:

[**C**L'21]:

[DKKLSS'19]:

mean estimation in time  $\tilde{O}(nd)/\text{poly}(\epsilon)$ .  $\tilde{O}(nd)$  time, heavy-tail mean estimation.

covariance estimation.

linear regression.

list-decodable mean estimation.

learning fixed-structure Bayesian networks.

robust gradient descent.

# Beyond Mean Estimation



# Outline

- Part I: Introduction
  - Motivation/Challenges
- Part II: Robust Algorithms for ML
  - High-Dimensional Statistics
  - Non-Convex Optimization
  - Learning with Strategic Agents



#### Part III: Future Directions

# Non-Convex Optimization

Extremely successful in practice.





# Non-Convex Optimization

Extremely successful in practice.

- In theory: NP-Hard.
- In practice: can be solved via (stochastic) gradient descent.

Why does non-convex optimization work?

- One possible explanation: all local optima are globally optimal.
- How robust are such landscape results?





# Matrix Completion

An unknown *n* by *n* matrix  $M^*$  with rank  $r \ll n$ . Input:  $M_{ij}^*$  for a set of observed entries  $(i, j) \in \Omega$ . Goal: Recover  $M^*$ .





# Matrix Completion

Application: recommendation systems (e.g., [Fiat et al. '01, Rennie and Srebro '05, Koren '09]):







# Matrix Completion

Application: recommendation systems

(e.g., [Fiat et al. '01, Rennie and Srebro '05, Koren '09]):

n users	0	-2 1	-2 0	-2 2	2-2	-2 1	-2 0	0	0	-2 1		1 0	0	-1 1	1	-1 -1	0	-1 1	-1 0	<i>r</i> features		
	1	0	-1	1	-1	0	-1	-1	1	0	L											
	2	0	-2	2	-2	0	-2	-2	 2	0												
	2	2	0	4	-4	2	0	-2	 2	2		We hope to recover $M^{\star}$										
	-1	1	2	0	0	1	2	1	-1	1	$\tilde{O}(nr)$ charactic											
	1	1	0	2	-2	1	0	-1	1	1	given $O(nr)$ obs								servations.			
	-2	2	4	0	0	2	4	2	-2	2												

*n* movies
### Previous Work

• Convex Relaxation (e.g., [Candès and Tao '10] [Recht '11])

• min 
$$||M||_{\star}$$
 s.t.  $M_{ij} = M_{ij}^{\star}$ 

• Runtime:  $O(n^4)$ .

In practice, people use non-convex approaches.

### Previous Work

• Non-Convex Approaches  $\min \|XX^{\top} - M^{\star}\|_{F}^{2}$ • min  $f(X) = \sum_{(i,j)\in\Omega} \left( (XX^{T})_{ij} - M_{ij}^{\star} \right)^{2}$  for  $X \in \mathbb{R}^{n \times r}$ .

If  $\widetilde{\Omega}(nr^6)$  entries are observed **uniformly at random**, any local optima X of f has  $XX^{\top} = M^*$  [Ge Lee Ma '16].



#### What if the observations are **not** uniformly at random?

## Semi-Random Adversary













## Robust Matrix Completion [C Ge '18]



Does not hurt **convex** approaches: min  $||M||_{\star}$  s.t.  $M_{ij} = M_{ij}^{\star}$ 

However, existing **non-convex** algorithms **are not robust** against such a semi-random adversary.



## Counter Examples [C Ge '18]

$$f(X) = \sum_{(i,j)\in\Omega} \left( M_{ij}^{\star} - (XX^{\top})_{ij} \right)^2$$
 has bad local optima.





-1

-1

1	1	1	1	-1	-1	-1	-1
1	1	1	1	-1	-1	-1	-1
1	1	1	1	-1	-1	-1	-1
1	1	1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	1	1	1
-1	-1	-1	-1	1	1	1	1
-1	-1	-1	-1	1	1	1	1
-1	-1	-1	-1	1	1	1	1

### Counter Examples [C Ge '18]



$$\begin{array}{ccc} 1 & & \begin{pmatrix} 1 \\ -1 \end{pmatrix} \rightarrow X = \begin{pmatrix} 1 \\ -(1-\epsilon) \end{pmatrix} \rightarrow \cdots \rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ & & & \\ \end{array}$$

$$\begin{array}{ccc} \bullet & & \\ -1 & & & \\ \end{array}$$

$$\begin{array}{ccc} X X^{\mathsf{T}} \approx \begin{pmatrix} 1 & -1+\epsilon \\ -1+\epsilon & 1-2\epsilon \end{pmatrix} \end{array}$$

$$f(X) = \sum_{(i,j)\in\Omega} \left( (XX^T)_{ij} - M^{\star}_{ij} \right)^2 \approx \|XX^T - M^{\star}\|_F^2$$

## Robust Matrix Completion [C Ge '18]

• Against a semi-random adversary:

Existing non-convex algorithms are not robust against such a semi-random adversary.

We can fix the non-convex approaches while preserving their efficiency.

# Robust Matrix Completion [C Ge '18]

If we know the adversary is changing the input like this ...





Preprocessing step: reweight the samples so that the input is "similar" to a random input.

$$f_{W}(X) = \sum_{(i,j)\in\Omega} W_{ij} (M_{ij}^{\star} - (XX^{\top})_{ij})^{2}$$

#### The graph formulation:





#### The graph formulation:



• Input: H = G(n, n, p) + extra edges





#### The graph formulation:



- Input: H = G(n, n, p) + extra edges
- Goal: Recover G

Recover a graph (spectrally) similar to G.

We show this can be solved in time:  $\tilde{O}(m \operatorname{poly}(1/\epsilon))$ .

 $(m = |E_H|, H' \approx_{\epsilon} G)$ 



Adapt algorithmic techniques developed for Graph Sparsification [Batson Spielman Srivastava '12, Lee and Sun '17].

## Robust Matrix Completion [C Ge '18]



# Robust Non-Convex Optimization

A meta framework for making non-convex approaches robust while preserving their efficiency.



## Outline

Part I: Introduction

• Motivation/Challenges

Part II: Robust Algorithms for ML

- High-Dimensional Statistics
- Non-Convex Optimization
- Learning with Strategic Agents



Part III: Future Directions

# Learning with Strategic Agents

- Data often come from agents, but traditionally we ignore the agents' incentives and consider learning problems in isolation.
- When the algorithms are used to make important decisions about the agents, agents may want to strategically provide data.
- Goal: design optimal algorithms or policies that take the agents' strategic behavior into consideration.

Motivating examples:

• Hiring committee



Motivating examples:

• Hiring committee



Motivating examples:

- Hiring committee
- College scout

The principal does not have direct access to the samples and relies on a strategic agent to provide samples.

What is the optimal policy?

Two types of agents: good



- Principal wants to accept the good type and reject the bad type.
- Agent wants to get accepted.

Principal has time to read m paper.

Each agent has *n* papers.

• Agent generates n samples from a publicly-known distribution based on his/her type, then choose which  $m \leq n$  samples to report.

## $Example \ 1 \ [Zhang \ \textbf{C} \ Conitzer \ '19]$

n = 50, m = 1. (b): top conference (b): average conference

0.05	0.95
0.005	0.995

Optimal policy: accept iff agent submits { ( )

- is accepted with prob.  $1 0.95^{50} \approx 0.92$ .
- is accepted with prob.  $1 0.995^{50} \approx 0.22$ .

# Example 2 [Zhang C Conitzer '19]









n = 3, m = 2.

Accept 
$$\{ \mathcal{B} \mathcal{B} \}$$
 and  $\{ \mathcal{B} \mathcal{B} \}!$ 

Not monotone in the likelihood ratio!

0.2

0.1

0.8

0.1

Yu Cheng

0

0.8

# Automated Mechanism Design

We study the structure and computation of the optimal policy.

- [Zhang C Conitzer '19]: Agents can withhold samples.
- [Zhang **C** Conitzer '19]:

Agents can transform the samples in limited ways.



# Automated Mechanism Design

We study the structure and computation of the optimal policy.

- [Zhang C Conitzer '19]: Agents can withhold samples.
- [Zhang **C** Conitzer '19]:

Agents can transform the samples in limited ways.

• [Zhang C Conitzer '21]:

Agents may have different preferences over outcomes. Generalization to more outcomes than accept/reject.

## Outline

- Part I: Introduction
  - Motivation/Challenges
- Part II: Robust Algorithms for ML
  - High-Dimensional Statistics
  - Non-Convex Optimization
  - Learning with Strategic Agents

#### Part III: Future Directions







0	-2	-2	-2	2	-2	-2	0
1	1	0	2	-2	1	0	-1
1	0	-1	1	-1	0	-1	-1
2	0	-2	2	-2	0	-2	-2
2	2	0	4	-4	2	0	-2
-1	1	2	0	0	1	2	1
1	1	0	2	-2	1	0	-1
-2	2	4	0	0	2	4	2















0	-2	-2	-2	2	-2	-2	0
1	1	0	2	-2	1	0	-1
1	0	-1	1	-1	0	-1	-1
2	0	-2	2	-2	0	-2	-2
2	2	0	4	-4	2	0	-2
-1	1	2	0	0	1	2	1
1	1	0	2	-2	1	0	-1
-2	2	4	0	0	2	4	2































Yu Cheng

# Future Direction I: Simple Algorithms

This talk:

- Polynomial time. X
- Fastest possible asymptotic runtime.  $\checkmark$

Future work: speed up the technology transfer into practice.

- Beyond fast asymptotic runtime.
- Leverage modern ML architectures (GPUs, DL packages).

## Future Direction I: Simple Algorithms

When a small fraction of the input is corrupted:

[Tukey '60, Huber '64]: provably robust statistical estimators.

[Diakonikolas+'16, Lai+'16]: polynomial-time robust statistical estimators.

[**C** Diakonikolas Ge '19]: provably robust estimators that are **as efficient as** their non-robust counterparts. Fast asymptotic runtime ≠ Practical Many algorithms are very sophisticated (e.g., ellipsoid method, SoS, JL lemma, matrix multiplication weight update) or require parameters that need careful tuning.

# High-Dimensional Robust Statistics

When a small fraction of the input is corrupted:

[Tukey '60, Huber '64]: provably robust statistical estimators.

[Diakonikolas+'16, Lai+'16]: polynomial-time robust statistical estimators.

[**C** Diakonikolas Ge '19]: provably robust estimators that are **as efficient as** their non-robust counterparts.

[**C** Diakonikolas Ge Soltanolkotabi '20]: provably robust estimators that can be computed using **standard first-order methods**. Robust Mean Estimation via Gradient Descent [C Diakonikolas Ge Soltanolkotabi '20]

- We consider non-convex formulations of robust mean estimation:
- Despite its non-convexity, we show that any (approximate) stationary point works!



```
for itr = 1:numItr
   Sigma_w_fun = @(v) X' * (w .* (X * v)) - (X' * w)^2 * v;
   [u, lambda] = eigs(Sigma_w_fun, d, 1);
   nabla_f_w = (X * u) .* (X * u) - 2 * (w' * (X * u)) * (X * u);
   w = w - stepSize * nabla_f_w / norm(nabla_f_w);
   w = project_onto_capped_simplex(w, 1 / (N - epsN));
end
```

Robust estimators that are **as efficient as** their non-robust counterparts.

[**C**DG'19]:

mean estimation in time  $\tilde{O}(nd)/\text{poly}(\epsilon)$ .

[DL'19][DHL'19]:

 $\tilde{O}(nd)$  time, heavy-tail mean estimation.

[CDGW'19][LY'20]: covariance estimation.

[C+'20]: linear regression. Span

Sparse mean?

[C+'20][D+'20]: list-decodable mean estimation.

[CL'21]: learning Bayesian networks. Robust estimators that can be computed using **standard first-order methods**.

[CDGS'20]: mean estimation via gradient descent. [HLZ'20]: heavy-tailed mean estimation. improved # of iterations via mirror descent. [Zhu+'21]: covariance estimation, linear regression. [CDKGGS'22, ongoing]: sparse mean estimation and sparse PCA.

Fast and simple?

### Future Direction II: Robust Deep Learning

This talk:

• Design and analysis of robust non-convex algorithms for problems that we understand very well w/o corruption.

Future work: robustness in deep learning.

- Expand the success of robust non-convex optimization.
- Develop heuristics based on theoretical analysis.

#### Future Direction II: Robust Deep Learning



[CGZ, ongoing]: one-bit and noisy matrix completion.[CDG, ongoing]: matrix sensing and phase retrieval.
## Future Direction II: Non-Convex Optimization

Robustness of neural networks:

- Under what assumptions can neural nets still converge to good parameters if 1% of the data is corrupted?
- If we want to throw away 1% of the training data, which 1% should we throw away?







## Future Directions III: General Theory of Robustness

- Systematic exploration of robustness.
  - Explore different adversarial models and examine whether commonly used algorithms are robust in these models.
  - Develop faster/simpler robust algorithms.
  - Translate the success of robust algorithm design to broader areas.
- Bridge the gap between the growing need for robust algorithms and the lack of general theory of robustness.









Ν

А

S

FΟ



U N I V E R S I T Y

1000

0.

IS ON

























**SC** 





































Part I: Introduction

• Motivation/Challenges

Part II: Robust Algorithms for ML

- High-Dimensional Statistics
- Non-Convex Optimization
- Learning with Strategic Agents

Part III: Future Directions

