# Reducing Certified Regression to Certified Classification for General Poisoning Attacks
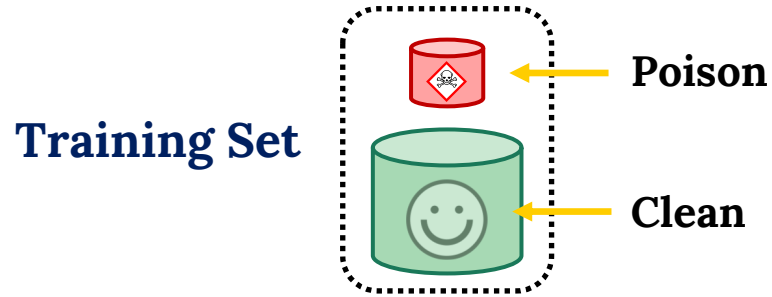
**Zayd Hammoudeh**    Daniel Lowd

SaTML 2023 – Raleigh, NC
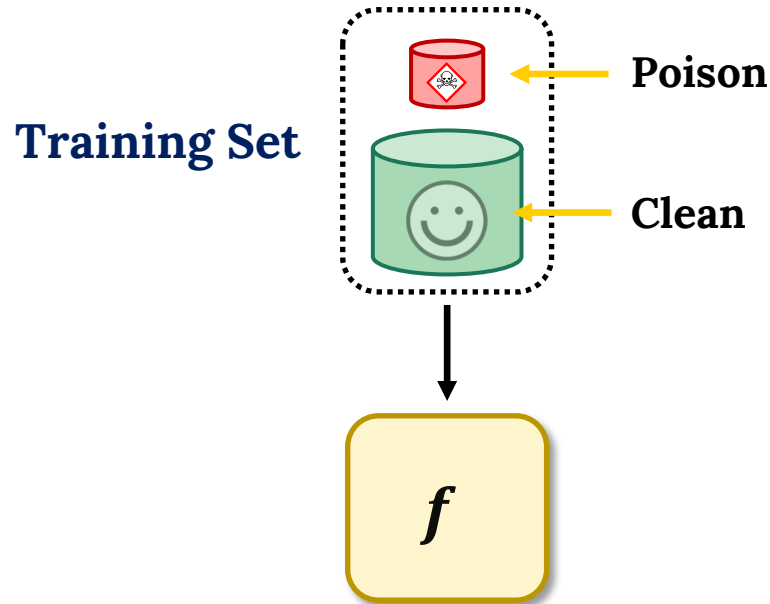
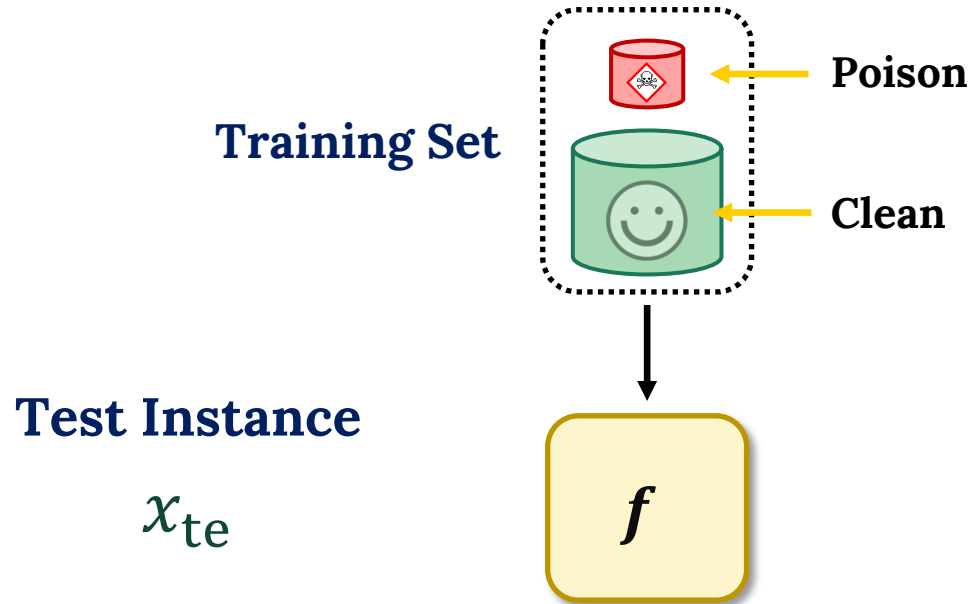# Data Poisoning Whirlwind Review

**Training Set**

Poison

Clean

# Data Poisoning Whirlwind Review

# Data Poisoning Whirlwind Review

# Data Poisoning Whirlwind Review



Training Set

Poison

Clean

Test Instance

$x_{\text{te}}$

$f$

Prediction

# Certified Regression against Poisoning

**Goal**: Certify *pointwise* robustness $R$ – the number of <u>arbitrary</u> instances that can be inserted or deleted from the training set with it guaranteed that:

$$\alpha \leq f(x_{\text{te}}) \leq \beta$$

- $\alpha, \beta \in \mathbb{R}$: User specified constants

# Certified Regression against Poisoning

**Goal**: Certify **_pointwise_** robustness $R$ – the number of <u>arbitrary</u> instances that can be inserted or deleted from the training set with it guaranteed that:

$$\alpha \leq f(x_{\text{te}}) \leq \beta$$

- $\alpha, \beta \in \mathbb{R}$: User specified constants

# **Structure of this Talk**

**Not the Focus**: Our six certified regressors

**Focus of this Talk**: Our reduction

# Specialized Robust Regressors Under Outliers & Poison

**Robust Linear ...**

**Robust ...**

**Efficient A...**

**High Dim...**

**Manipulating ...
and Counter...**

Matthew Jagielski*, Alina ...

*Northeastern University, Boston, ...

arXiv:1803.03241v3 [cs.LG] 4 Jun 2020

arXiv:1805.11643v3 [cs.LG] 29 May 2019

arXiv:1804.00308v3 [cs.CR] 28 Sep 2021

arXiv:1301.2725v1 [stat.ML] 12 Jan 2013

## Robust High Dimensional Sparse Regression and Matching Pursuit

Yudong Chen, Constantine Caramanis and Shie Mannor

### Abstract

In this paper we consider high dimensional sparse regression, and develop strategies able to deal with arbitrary – possibly, severe or coordinated – errors in the covariance matrix $X$. These may come from corrupted data, persistent experimental errors, or malicious respondents in surveys/recommender systems, etc. Such non-stochastic error-in-variables problems are notoriously difficult to treat, and as we demonstrate, the problem is particularly pronounced in high-dimensional settings where the primary goal is *support recovery* of the sparse regressor. We develop algorithms for support recovery in sparse regression, when some number $n_1$ out of $n + n_1$ total covariate/response pairs are *arbitrarily (possibly maliciously) corrupted*. We are interested in understanding how many outliers, $n_1$, we can tolerate, while identifying the correct support. To the best of our knowledge, neither standard outlier rejection techniques, nor recently developed robust regression algorithms (that focus only on corrupted response variables), nor recent algorithms for dealing with stochastic noise or erasures, can provide guarantees on support recovery. Perhaps surprisingly, we also show that the natural brute force algorithm that searches over all subsets of $n$ covariate/response pairs, and all subsets of possible support coordinates in order to minimize regression error, is remarkably poor, unable to correctly identify the support with even $n_1 = O(n/(\sqrt{k}\log p))$ corrupted points, where $k$ is the sparsity. This is true even in the basic setting we consider, where all authentic measurements and noise are independent and sub-Gaussian. In this setting, we provide a simple algorithm – no more computationally taxing than OMP – that gives stronger performance guarantees, recovering the support with up to $n_1 = O(n/(\sqrt{k}\log p))$ corrupted points, where $p$ is the dimension of the signal to be recovered.

### I. INTRODUCTION

Linear regression and sparse linear regression seek to express a response variable as the linear combination of (a small number of) covariates. They form one of the most basic procedures in statistics, engineering, and science. More recently, regression has found increasing applications in the high-dimensional regime, where the number of variables, $p$, is much larger than the number of measurements or observations, $n$. Applications in biology, genetics, as well as in social networks, human behavior prediction and recommendation, abound, to name just a few. The key structural property exploited in high-dimensional regression, is that the regressor is often sparse, or near sparse, and as much recent research has demonstrated, in many cases it can be efficiently recovered, despite the grossly underdetermined nature of the problem (e.g., [8], [6], [4], [12], [31]). Another common theme in large-scale learning problems – particularly problems in the high-dimensional regime – is that we not only have big data, but we have dirty data. Recently, attention has focused on the setting where the output (or response) variable and the matrix of covariates are plagued by erasures, and/or by stochastic additive noise [23], [26], [27], [9], [10]. Yet many applications, including those mentioned, may suffer from persistent errors, that are ill-modeled by stochastic distribution; indeed, many applications, particularly those modeling human behavior, may exhibit maliciously corrupted data.

This paper is about extending the power of regression, and in particular, sparse high-dimensional regression, to be robust to this type of noise. We call this *deterministic* or *cardinality constrained* robustness, because rather than restricting the magnitude of the noise, or any other such property of the noise, we merely assume there is some bounded number of data points, or how many covariates, are corrupted. Other than this number, we make absolutely no assumptions on what the adversary can do – the adversary is virtually unlimited in computational power and knowledge

5

# Adversarially Robust Regressors Make **Strong Assumptions**

## Data Distribution Assumptions

- Sparsity/low rank
- Linear data distribution with AWGN

## Model architecture assumptions

- Linear model

## Distributional Guarantees Only

- No insight into individual predictions' robustness

# Our Goal

Provably robust regressors that are **general**:

◉ No data distribution assumptions

◉ Model architecture agnostic

◉ Stop reinventing the wheel.
  ○ **Consistently** state–of–the–art with **minimal effort**

# A Bit of a Detour

# Certified Poisoning **Classifiers** Show Promise

# Certified Poisoning **Classifiers** Show Promise

# **Strengths of Certified Poisoning <span style="color:red">Classifiers</span>**

1. No data distribution assumptions

2. Model architecture agnostic

3. Strong empirical performance
   - Certify 65% of MNIST predictions up to 0.8% arbitrary poison

   - Certify 16% of CIFAR10 predictions up to 0.1% arbitrary poison

# General Structure of a Certified Poisoning Classifier

# General Structure of a Certified Poisoning Classifier

**Certified Classifier** $(f)$



**Test Instance** $x_{\text{te}}$

**Votes Generator**

0
0
0
1
0
1
0

**Plurality Label**

**Prediction** $f(x_{\text{te}})$

**Robustness Certifier**

$R$

**Certified Robustness**

**Multiset of Votes** $(\mathcal{V})$

# General Structure of a Certified Poisoning Classifier



**Certified Classifier** ($f$)

**Test Instance** $x_{\text{te}}$

**Votes Generator**

**Plurality Label**

**Robustness Certifier**

**Prediction** $f(x_{\text{te}})$

$R$

**Certified Robustness**

**Multiset of Votes** ($\mathcal{V}$)

# General Structure of a Certified Poisoning Classifier



Certified Classifier ($f$)

Test Instance $x_{\text{te}}$

Votes Generator

Plurality Label

Prediction $f(x_{\text{te}})$

Robustness Certifier

$R$

Certified Robustness

Multiset of Votes ($\mathcal{V}$)

# General Structure of a Certified Poisoning Classifier



Certified Classifier ($f$)

Test Instance $x_{\text{te}}$

Votes Generator

Plurality Label

Prediction $f(x_{\text{te}})$

Robustness Certifier

$R$

Certified Robustness

Multiset of Votes ($\mathcal{V}$)

# General Structure of a Certified Poisoning Classifier



Certified Classifier $(f)$

Test Instance $x_{\text{te}}$

Votes Generator

Plurality Label

Prediction $f(x_{\text{te}})$

Robustness Certifier

$R$

Certified Robustness

Multiset of Votes $(\mathcal{V})$

# General Structure of a Certified Poisoning Classifier

**Certified Classifier** $(f)$

**Test Instance** $x_{\text{te}}$

**Votes Generator**

**Plurality Label**

**Prediction** $f(x_{\text{te}})$

**Robustness Certifier**

$R$

**Certified Robustness**

**Multiset of Votes** $(\mathcal{V})$

# General Structure of a Certified Poisoning Classifier

# General Structure of a Certified Poisoning Classifier



Certified Classifier ($f$)

Test Instance $x_{\text{te}}$

Votes Generator

Plurality Label

Prediction $f(x_{\text{te}})$

Robustness Certifier

$R$

Certified Robustness

Multiset of Votes ($\mathcal{V}$)

# **Example:** kNN Certified Classifier [Jia+22]

# **Example:** kNN Certified Classifier [Jia+22]



**Multiset of Votes** $(\mathcal{V})$

$x_{\text{te}}$

$k = 5$

# **Example:** kNN Certified Classifier [Jia+22]



**Multiset of Votes** $(\mathcal{V})$

$x_{\text{te}}$

$k = 5$

# **Example:** kNN Certified Classifier [Jia+22]

$$0 \quad 0 \quad 0 \quad 0 \quad 1$$

**Vote Distribution:**

◉ 4 votes label $0$

◉ 1 vote label $1$

**Robustness Certifier:** At least two votes much change to perturb the plurality label

$$R = \left\lfloor \frac{4 - 1}{2} \right\rfloor = 1$$

# Example: Certified Ensemble Classifier [LF21]

# **Example:** Certified Ensemble Classifier [LF21]

**Ensemble**

$f_1$

$f_2$

$f_3$

$f_4$

$f_5$

# **Example:** Certified Ensemble Classifier [LF21]

**Test Instance**     **Ensemble**

$x_\text{te}$     $\boxed{f_1}$

$x_\text{te}$     $\boxed{f_2}$

$x_\text{te}$     $\boxed{f_3}$

$x_\text{te}$     $\boxed{f_4}$

$x_\text{te}$     $\boxed{f_5}$

# **Example:** Certified Ensemble Classifier [LF21]

# **Example:** Certified Ensemble Classifier [LF21]



**Test Instance**   **Ensemble**

$x_{\text{te}}$ → $f_1$ → 0

$x_{\text{te}}$ → $f_2$ → 0

$x_{\text{te}}$ → $f_3$ → 0

$x_{\text{te}}$ → $f_4$ → 1

$x_{\text{te}}$ → $f_5$ → 0

**Votes Multiset** $(\mathcal{V})$

**Robustness Certification**

15

# **Example:** Certified Ensemble Classifier [LF21]



**Test Instance**  **Ensemble**

$x_\text{te} \rightarrow \boxed{f_1} \rightarrow 0$

$x_\text{te} \rightarrow \boxed{f_2} \rightarrow 0$

$x_\text{te} \rightarrow \boxed{f_3} \rightarrow 0$

$x_\text{te} \rightarrow \boxed{f_4} \rightarrow 1$

$x_\text{te} \rightarrow \boxed{f_5} \rightarrow 0$

**Votes Multiset** $(\mathcal{V})$

**Robustness Certification**

**Vote Distribution:**
- 4 votes label  0
- 1 vote label  1

15

# **Example:** Certified Ensemble Classifier [LF21]



**Test Instance**   **Ensemble**

$x_{\text{te}} \longrightarrow f_1 \longrightarrow 0$

$x_{\text{te}} \longrightarrow f_2 \longrightarrow 0$

$x_{\text{te}} \longrightarrow f_3 \longrightarrow 0$

$x_{\text{te}} \longrightarrow f_4 \longrightarrow 1$

$x_{\text{te}} \longrightarrow f_5 \longrightarrow 0$

**Votes Multiset** $(\mathcal{V})$

**Robustness Certification**

**Vote Distribution**:
- 4 votes label $0$
- 1 vote label $1$

**Certified Robustness**:
$$R = \left\lfloor \frac{4-1}{2} \right\rfloor = 1$$

15

# Reducing Certified Regression to Voting-Based Certified Classification

# "Don't Reinvent the Wheel"

**Reduction**: An algorithm for converting a problem $Q$ into a **different problem** $Q'$ that can be **readily solved**.

) to certified classification ($Q'$)

*"Transform certified poisoning classifiers into certified regressors"*

**Benefits**:
- Inherit the strengths of the certified classifiers
- Each improved certified classifier improves certified regression

# "Don't Reinvent the Wheel"

)

**Reduction**: An algorithm for converting a problem $Q$ into a **different problem** $Q'$ that can be **readily solved**.

**Our Idea:** Reduce certified regression ($Q$) to certified classification ($Q'$)

*"Transform certified poisoning classifiers into certified regressors"*

*"Transform certified poisoning classifiers into certified regressors"*

**Benefits**:

- Inherit the strengths of the certified classifiers
- Each improved certified classifier improves certified regression

# "Don't Reinvent the Wheel"

)

**Reduction**: An algorithm for converting a problem $Q$ into a **different problem** $Q'$ that can be **readily solved**.

**Our Idea:** Reduce certified regression ($Q$) to certified classification ($Q'$)

*"Transform certified poisoning classifiers into certified regressors"*

*"Transform certified poisoning classifiers into certified regressors"*

**Benefits**:

- Inherit the strengths of the certified classifiers
- Each improved certified classifier improves certified regressionEach improved certified classifier improves certified

# Key Insight of the Reduction

For any binary multiset,
the **plurality label** and **median** have
**equivalent robustness**

# Relating Real-Valued and Binary Robustness

**Multiset of Real Votes** $(\mathcal{V})$    2.2    3.1    4.2    5.0    6.1

# Relating Real-Valued and Binary Robustness

**Multiset of Real Votes** $(\mathcal{V})$   $\boxed{2.2}$  $\boxed{3.1}$  $\boxed{4.2}$  $\boxed{5.0}$  $\boxed{6.1}$

med $\mathcal{V}$

# Relating Real-Valued and Binary Robustness

**Upper Bound**

$\beta$

**Multiset of
Real Votes** $(\mathcal{V})$   2.2   3.1   4.2   5.0   6.1

med $\mathcal{V}$

# Relating Real-Valued and Binary Robustness

**Upper Bound**

$\beta$

**Multiset of Real Votes** $(\mathcal{V})$    2.2    3.1    4.2    5.0    6.1

med $\mathcal{V}$

# Relating Real-Valued and Binary Robustness

**Upper Bound**

$\beta$

**Robustness Certifier**

**Multiset of Real Votes** ($\mathcal{V}$)

2.2   3.1   4.2   5.0   6.1

med $\mathcal{V}$

# Relating Real-Valued and Binary Robustness

**Upper Bound**

$\beta$

**Multiset of Real Votes** ($\mathcal{V}$)

| 2.2 | 3.1 | 4.2 | 5.0 | 6.1 |

med $\mathcal{V}$

**Vote Distribution**:

⬤ 4 votes label 0

⬤ 1 vote label 1

# Relating Real-Valued and Binary Robustness

**Upper Bound**

$\beta$

**Multiset of Real Votes** ($\mathcal{V}$)

2.2   3.1   4.2   5.0   6.1

med $\mathcal{V}$

**Robustness Certifier**

**Vote Distribution**:

⊙  4 votes label  0

⊙  1 vote label  1

**Certified Robustness**:

$$R = \left\lfloor \frac{4-1}{2} \right\rfloor = 1$$

"Transform certified poisoning classifiers into certified regressors"

# General Structure of a Certified Poisoning Classifier



Certified Classifier ($f$)

Test Instance $x_{\text{te}}$

Votes Generator

Multiset of Votes ($\mathcal{V}$)

Plurality Label

Robustness Certifier

Prediction $f(x_{\text{te}})$

$R$

Certified Robustness

# General Structure of a Certified Poisoning **Regressor**

**Certified Regressor** $(f)$

# General Structure of a Certified Poisoning Regressor



Certified Regressor ($f$)

Test Instance $x_{te}$

Votes Generator

Multiset of Votes ($\mathcal{V}$)

0
0
0
1
0

Plurality Label → Prediction $f(x_{te})$

Robustness Certifier → $R$ Certified Robustness

# General Structure of a Certified Poisoning Regressor



Certified Regressor ($f$)

Test Instance $x_{\text{te}}$ → Votes Generator → [2.2, 3.1, 4.2, 5.0, 6.1]

Plurality Label → Prediction $f(x_{\text{te}})$

Robustness Certifier → $R$ Certified Robustness

Real Multiset of Votes ($\mathcal{V}$)

# General Structure of a Certified Poisoning **Regressor**



Certified **Regressor** ($f$)

Test Instance $x_{\text{te}}$

Votes Generator

2.2
3.1
4.2
5.0
6.1

**Median** med $\mathcal{V}$

Prediction $f(x_{\text{te}})$

**Robustness Certifier**

$R$

Certified Robustness

**Real** Multiset of Votes ($\mathcal{V}$)

# General Structure of a Certified Poisoning Regressor

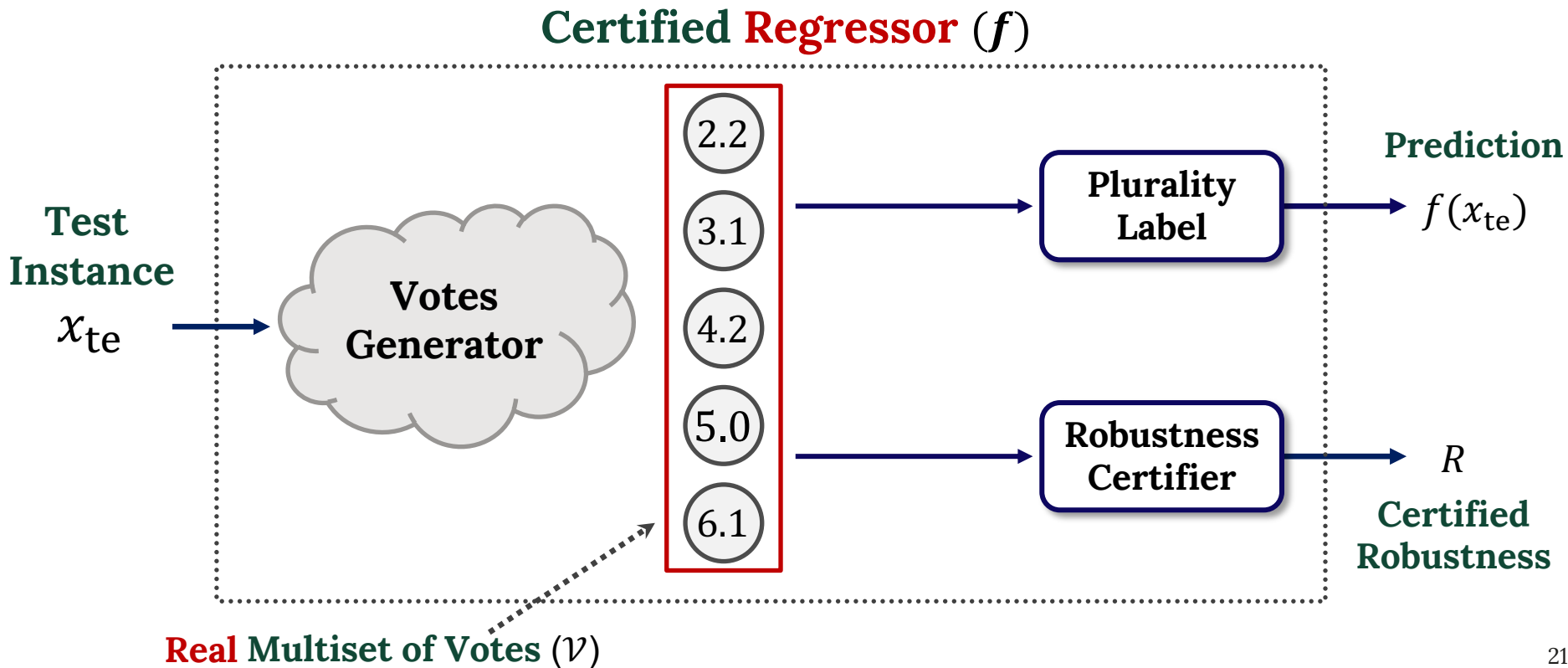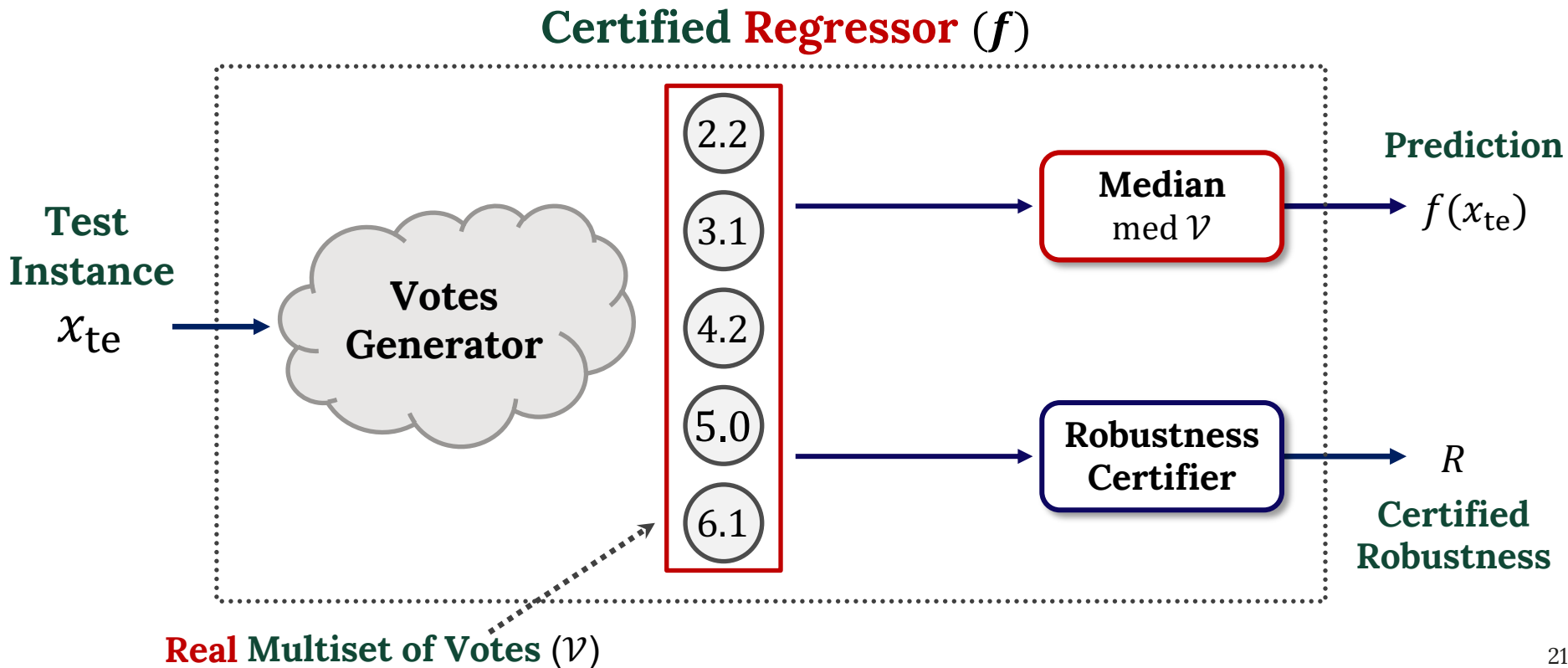# General Structure of a Certified Poisoning Regressor



Certified Regressor ($f$)

Test Instance $x_{\text{te}}$

Votes Generator

2.2
3.1
4.2
5.0
6.1

Median med $\mathcal{V}$

Prediction $f(x_{\text{te}})$

$\leq \beta$

Robustness Certifier

$R$

Certified Robustness

Real Multiset of Votes ($\mathcal{V}$)

Binarize $\mathcal{V}$

# **Summary: Reducing Certified Regression to Voting-Based Certified Classification**

## **Three simple steps**:

- Generate real-valued votes instead of labels

- Use median as the decision function

- Binarize the real-valued votes $\mathcal{V}$ using threshold $\beta$

# Our Reduction Yields a Suite of Certified Regressors

We propose six certified regressors:

- Two based on certified **nearest neighbor** classifiers [Jia+22]

- Two based on certified **ensemble** classifiers [LF21, WLF22]

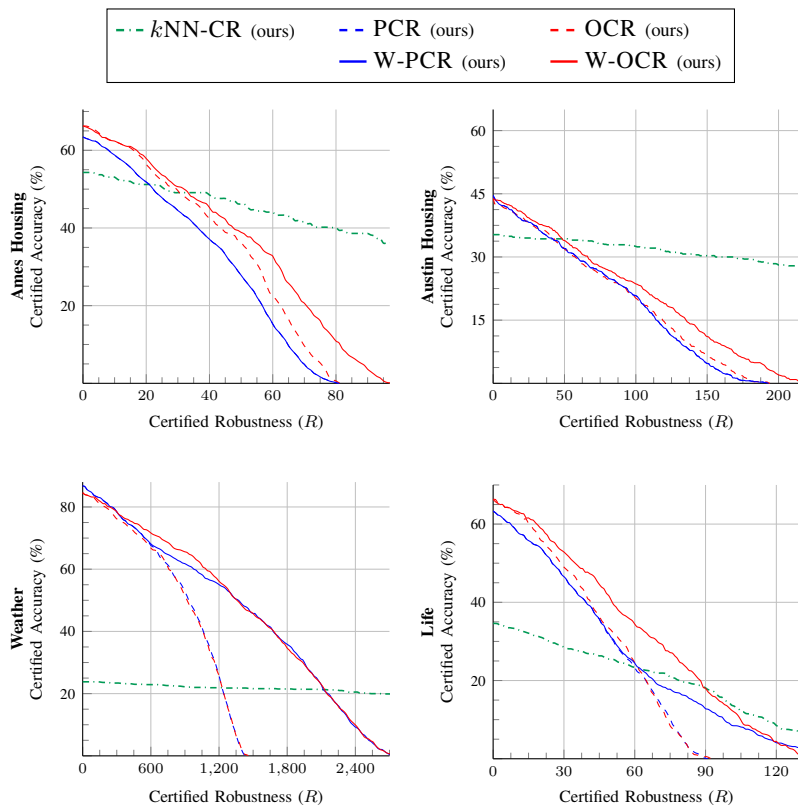- Two based on **our improved** certified ensemble classifiers

# Empirical Evaluation

# Empirical Evaluation

⦿ **Datasets**: 5 Regression + 1 Binary Classification

⦿ **Performance Metric**: **Certified accuracy**
  ○ Percentage of correctly predicted test instances given $\alpha$ and $\beta$ with certified robustness $R \geq \psi$

⦿ **Model Architecture Agnostic:** Decision trees and linear models

# Certified Regression – Takeaways



**Certified Accuracy**:
- Half of predictions up to 1% poisoning

- Third of predictions up to 4% poisoning

**Method Comparison**:
- **Nearest Neighbors**: Better maximum robustness ($R$)

- **Ensemble**: Better accuracy

# One more thing...

# Reducing Certified Regression to Certified Classification for General Poisoning Attacks

**Zayd Hammoudeh**    Daniel Lowd

SaTML 2023 – Raleigh, NC

UNIVERSITY OF
OREGON

# Reducing Certified Regression to Certified Classification ~~for General Poisoning Attacks~~

**Zayd Hammoudeh**    Daniel Lowd

SaTML 2023 – Raleigh, NC

UNIVERSITY OF
OREGON

# References

[LF21] A. Levine and S. Feizi, "Deep partition aggregation: Provable defenses against general poisoning attacks," ICLR, 2021.

[Jia+22] J. Jia, Y. Liu, X. Cao, and N. Gong. "Certified Robustness of Nearest Neighbors against Data Poisoning and Backdoor Attacks" AAAI, 2022.

[WLF22] W. Wang, A. Levine, and S. Feizi "Improved Certified Defenses against Data Poisoning with (Deterministic) Finite Aggregation" ICML, 2022.