

Measuring the Effect of Training Data on Deep Learning Predictions via Randomized Experiments

Jinkun Lin (NYU)*, Anqi Zhang (NYU)*,

Mathias Lécuyer (UBC), Jinyang Li (NYU), Aurojit Panda (NYU), Siddhartha Sen (MSR)



NEW YORK UNIVERSITY

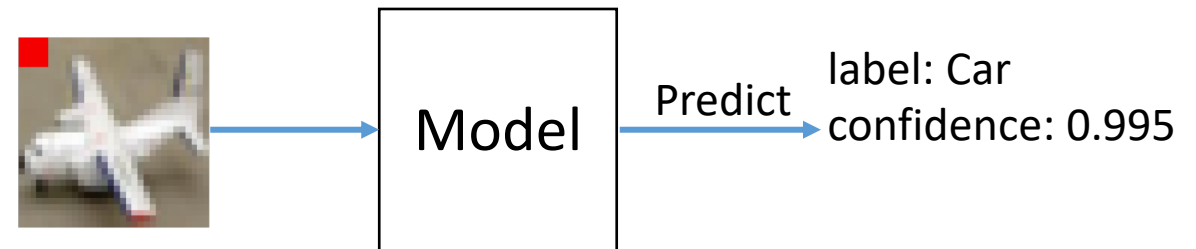


THE UNIVERSITY
OF BRITISH COLUMBIA

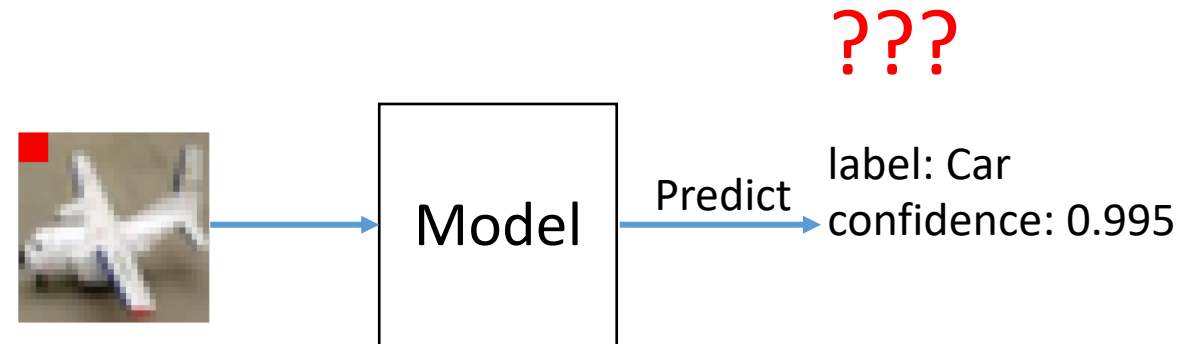
Microsoft®

Research

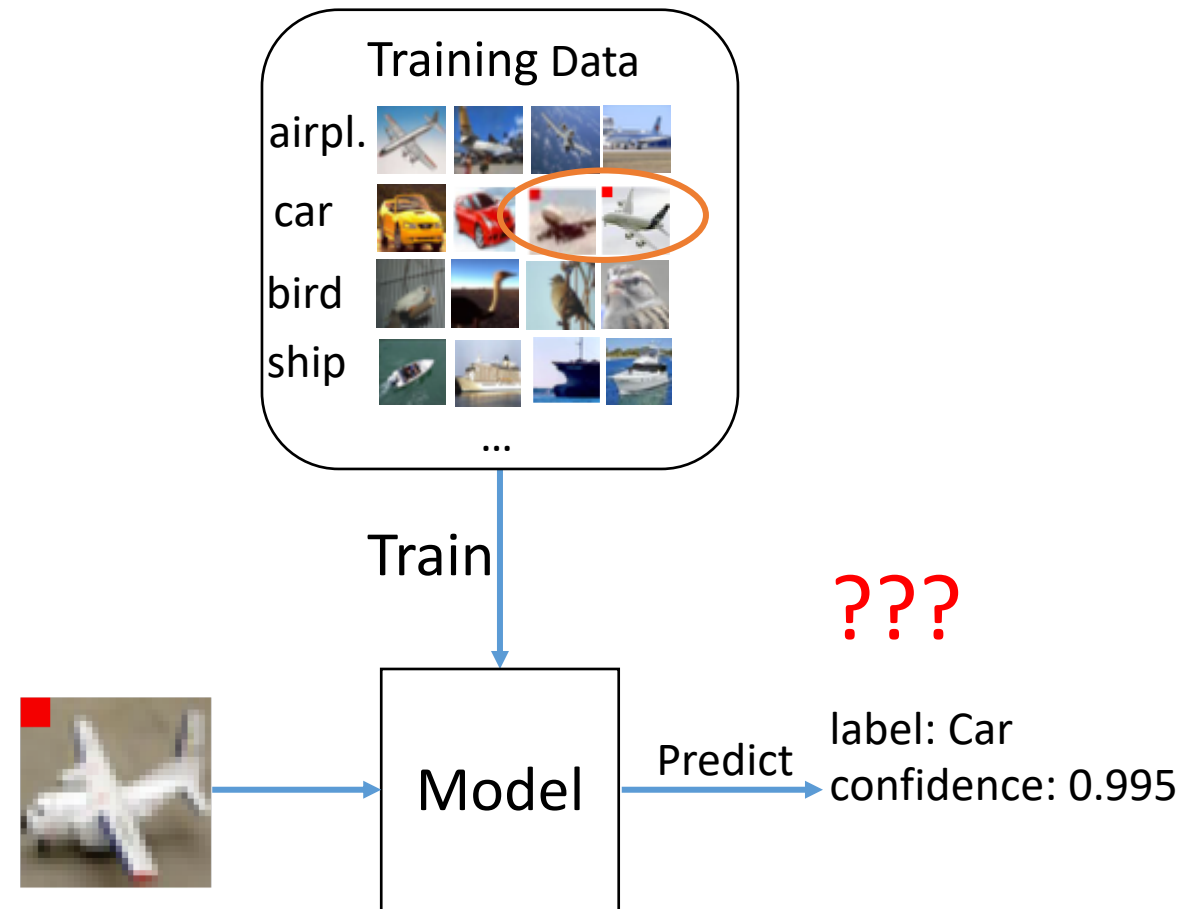
Problem statement



Problem statement

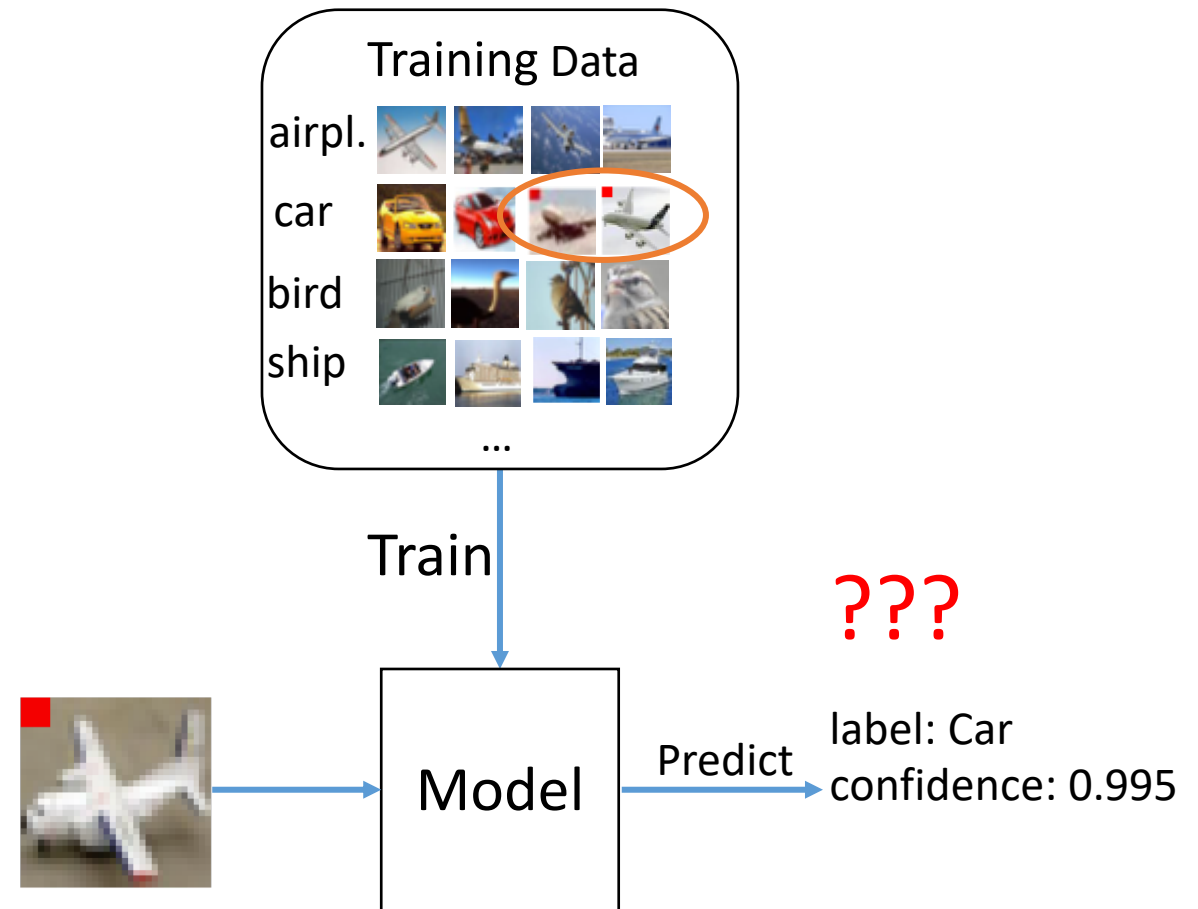


Problem statement



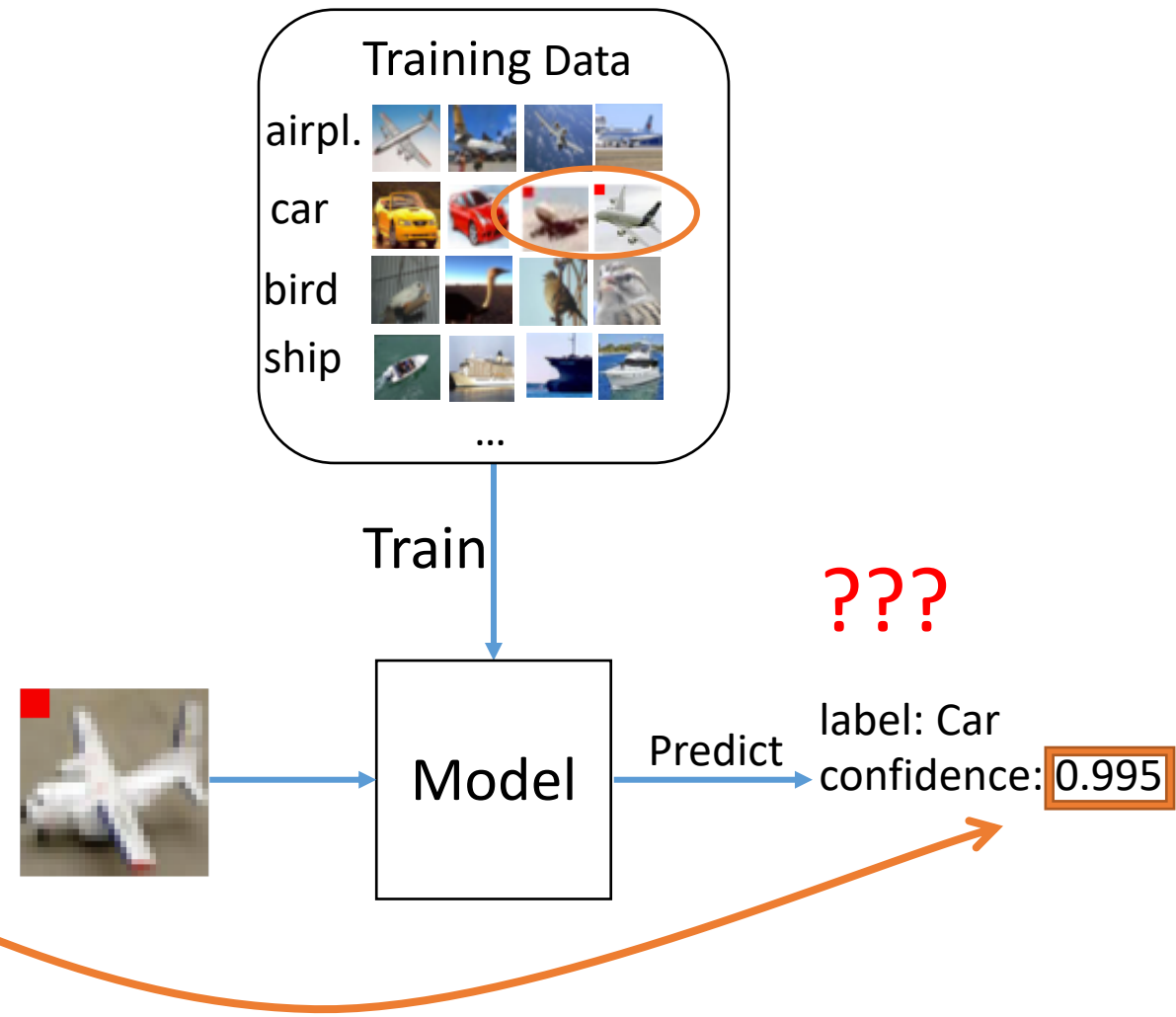
Problem statement

- How do we pinpoint the training examples that contribute significantly to a behavior?



Problem statement

- How do we pinpoint the training examples that contribute significantly to a behavior?
- Behavior is usually quantified by some utility:
 - E.g., $U(\cdot)$ = confidence score for a prediction



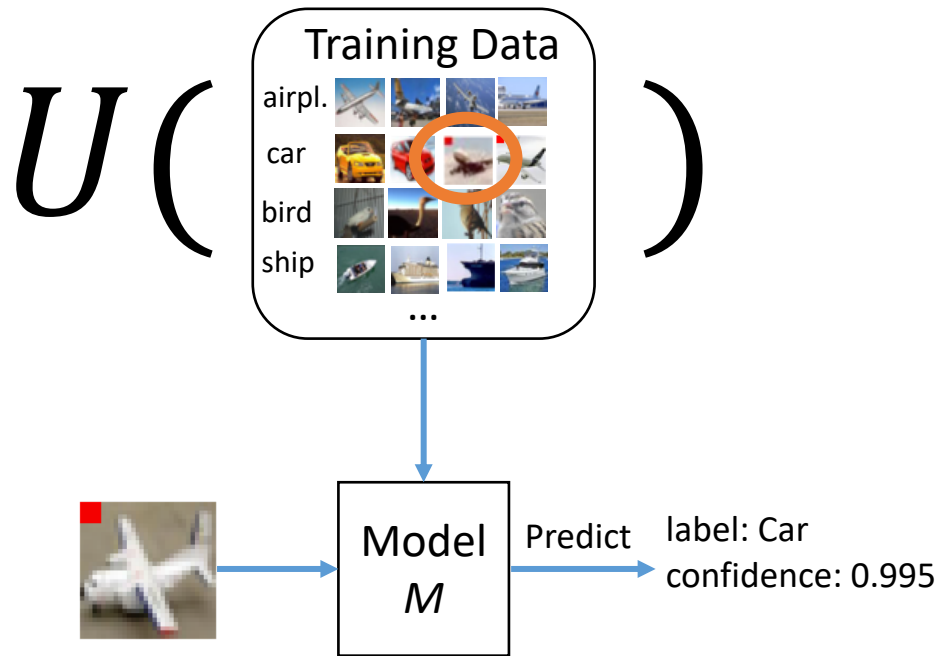
Influence Functions

- Influence function: measures the contribution of  using marginal effect



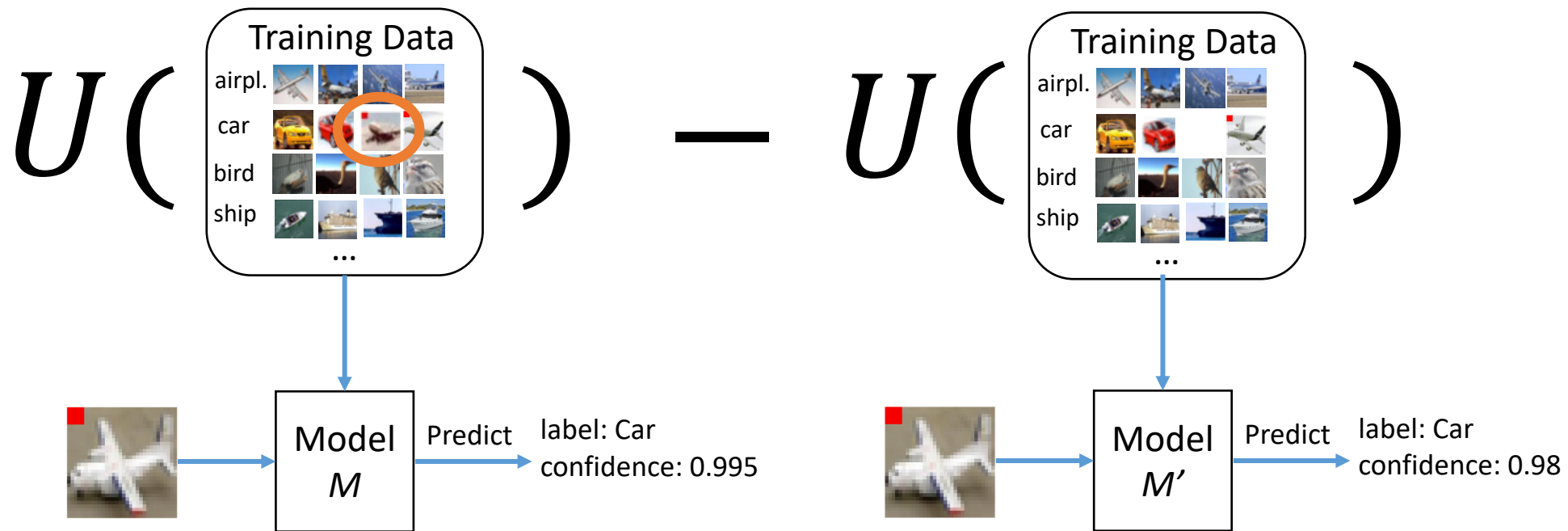
Influence Functions

- Influence function: measures the contribution of  using marginal effect



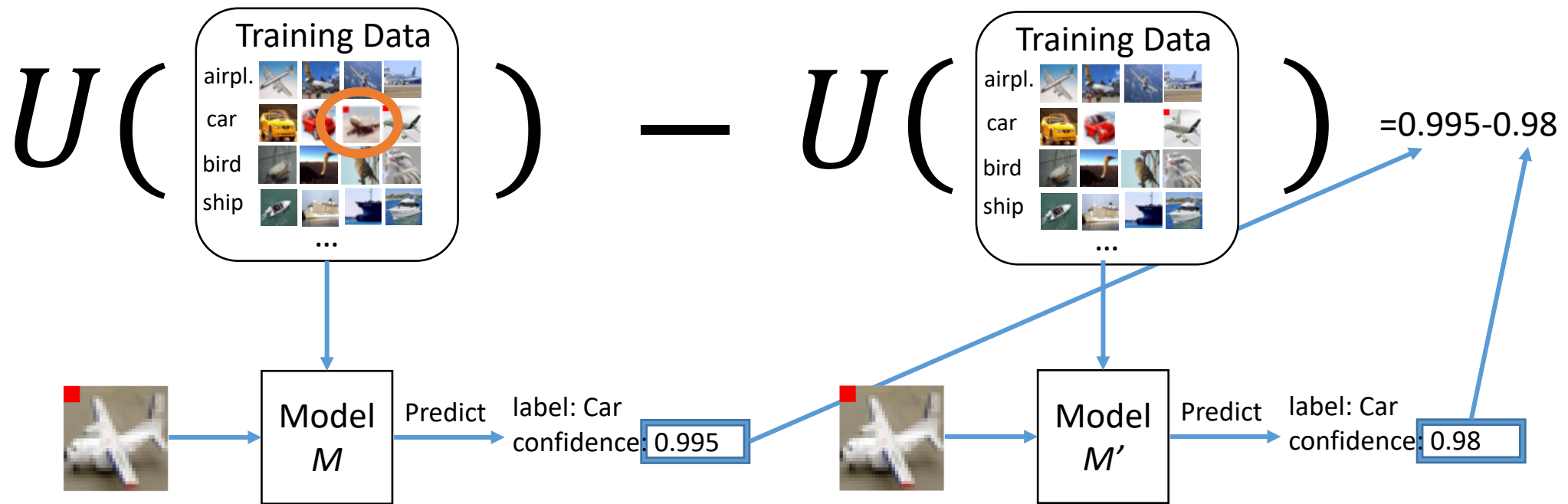
Influence Functions

- Influence function: measures the contribution of  using marginal effect



Influence Functions

- Influence function: measures the contribution of  using marginal effect



Why influence functions fall short?

- Influence function: measures the contribution of  using marginal effect

$$U\left(\begin{array}{c} \text{Training Data} \\ \text{airpl.} \begin{array}{cccc} \text{img1} & \text{img2} & \text{img3} & \text{img4} \end{array} \\ \text{car} \begin{array}{cccc} \text{img5} & \text{img6} & \text{img7} & \text{img8} \end{array} \\ \text{bird} \begin{array}{cccc} \text{img9} & \text{img10} & \text{img11} & \text{img12} \end{array} \\ \text{ship} \begin{array}{cccc} \text{img13} & \text{img14} & \text{img15} & \text{img16} \end{array} \\ \dots \end{array}\right) - U\left(\begin{array}{c} \text{Training Data} \\ \text{airpl.} \begin{array}{cccc} \text{img1} & \text{img2} & \text{img3} & \text{img4} \end{array} \\ \text{car} \begin{array}{cccc} \text{img5} & \text{img6} & \text{img7} & \text{img8} \end{array} \\ \text{bird} \begin{array}{cccc} \text{img9} & \text{img10} & \text{img11} & \text{img12} \end{array} \\ \text{ship} \begin{array}{cccc} \text{img13} & \text{img14} & \text{img15} & \text{img16} \end{array} \\ \dots \end{array}\right) = 0.995 - 0.98$$

- 👍 Efficient: no model training needed assuming convexity

Why influence functions fall short?

- 🙅 DNNs are not convex
- 🙅 Marginal effect is *close to 0* regardless of the data point

$$|U(\text{img1}) - U(\text{img2})| \approx 0$$

Why influence functions fall short?

- 🙌 DNNs are not convex
- 🙌 Marginal effect is *close to 0* regardless of the data point

$$|U(\text{grid of 16 images}) - U(\text{grid of 16 images})| \approx 0$$



Observation: marginal effect is *prominent* when removing from **smaller subsets** of the training set!

$$|U(\text{grid of 4 images}) - U(\text{grid of 4 images})| > 0$$

Average Marginal Effect (AME)

- Sample various subsets,

$$S \sim \{ \begin{array}{|c|} \hline \text{[Grid of 16 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 10 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \dots \}$$

Average Marginal Effect (AME)

- Sample various subsets, Average their Marginal Effects (AME)

$$E_S[U(S + \{i\}) - U(S)], S \sim \{ \begin{array}{|c|} \hline \text{[Grid of 16 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 10 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 5 small images]} \\ \hline \end{array} \quad \dots \}$$

Average Marginal Effect (AME)

- Sample various subsets, Average their Marginal Effects (AME)

$$E_S[U(S + \{i\}) - U(S)], S \sim \{ \begin{array}{|c|} \hline \text{img1} \\ \hline \end{array} \begin{array}{|c|} \hline \text{img2} \\ \hline \end{array} \begin{array}{|c|} \hline \text{img3} \\ \hline \end{array} \begin{array}{|c|} \hline \text{img4} \\ \hline \end{array} \dots \}$$

- Each subset is drawn by including each data point independently with probability $p \sim \mathcal{P}$

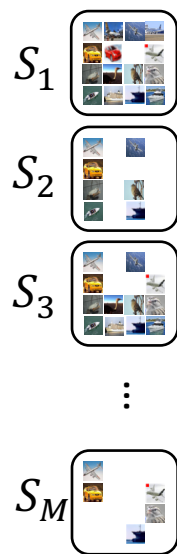
Average Marginal Effect (AME)

- Sample various subsets, Average their Marginal Effects (AME)

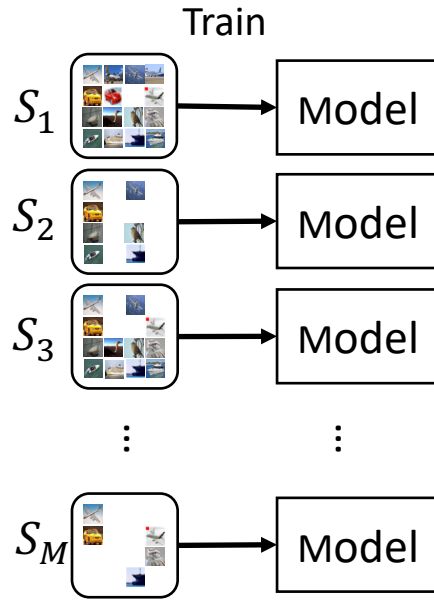
$$E_S[U(S + \{i\}) - U(S)], S \sim \{ \begin{array}{|c|} \hline \text{img1} \\ \hline \end{array} \begin{array}{|c|} \hline \text{img2} \\ \hline \end{array} \begin{array}{|c|} \hline \text{img3} \\ \hline \end{array} \begin{array}{|c|} \hline \text{img4} \\ \hline \end{array} \dots \}$$

- Each subset is drawn by including each data point independently with probability $p \sim \mathcal{P}$
- Estimate via randomized experiments and LASSO regression

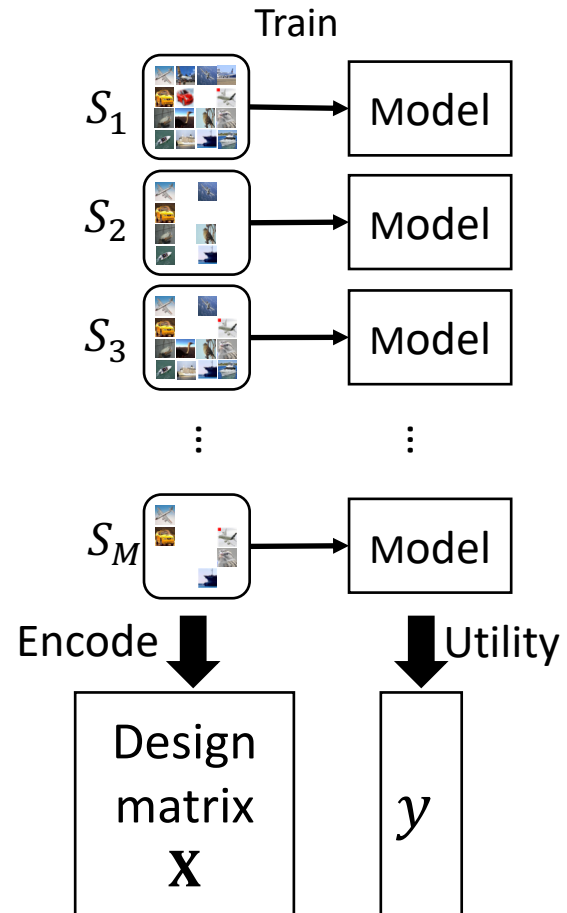
Efficient Estimation of AME



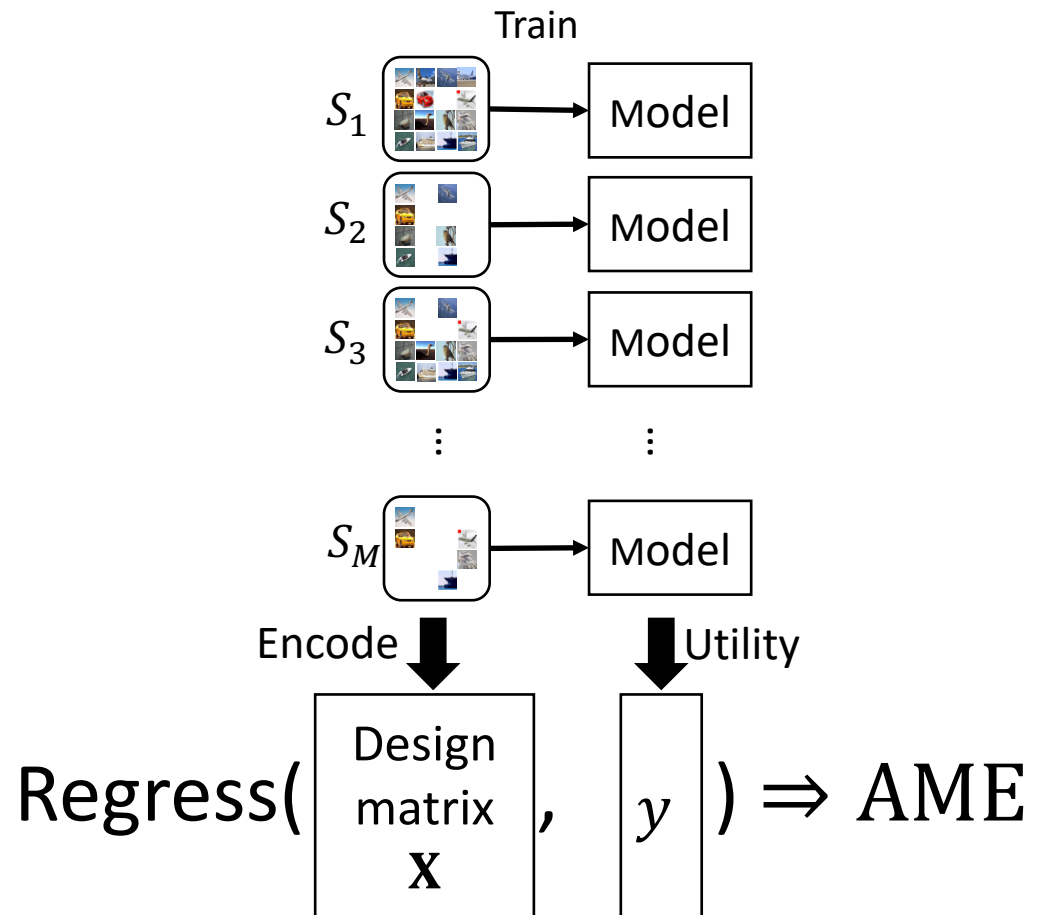
Efficient Estimation of AME



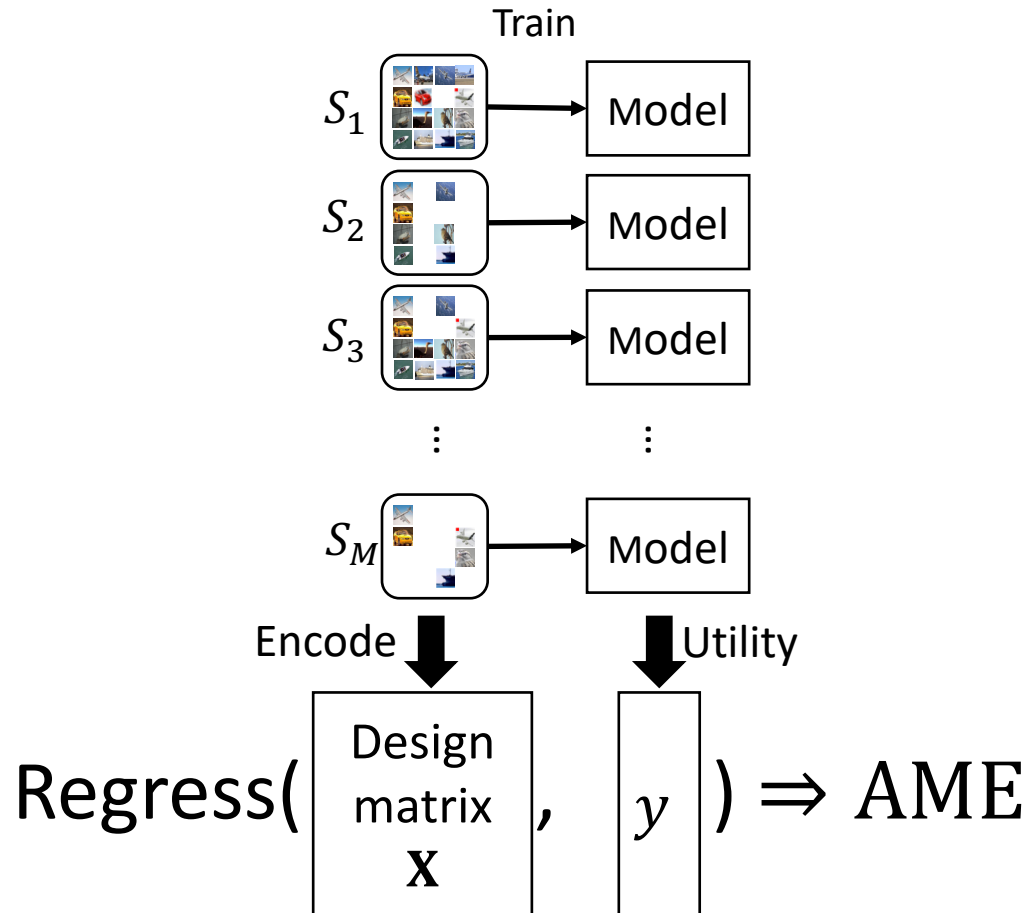
Efficient Estimation of AME



Efficient Estimation of AME



Efficient Estimation of AME



- **Sparsity assumption:** a few (k) data points have large contribution
- **Scalable:** train $O(k \log N)$ models where N : training set size

Sampling Distribution \mathcal{P} for AME

$$E_S[U(S + \{i\}) - U(S)], S \sim \{ \begin{array}{|c|} \hline \text{[Grid of 16 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 6 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 10 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \dots \}$$

- Subsets are drawn by including each data point independently with probability $p \sim \mathcal{P}$

Sampling Distribution \mathcal{P} for AME

$$E_S[U(S + \{i\}) - U(S)], S \sim \{ \begin{array}{|c|} \hline \text{[Grid of 16 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 9 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \dots \}$$

- Subsets are drawn by including each data point independently with probability $p \sim \mathcal{P}$
- $\mathcal{P} = \text{Uniform}(0,1)$, AME = SV, but incompatible with fast estimation

Sampling Distribution \mathcal{P} for AME

$$E_S[U(S + \{i\}) - U(S)], S \sim \{ \begin{array}{|c|} \hline \text{[Grid of 16 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 9 small images]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Grid of 4 small images]} \\ \hline \end{array} \quad \dots \}$$

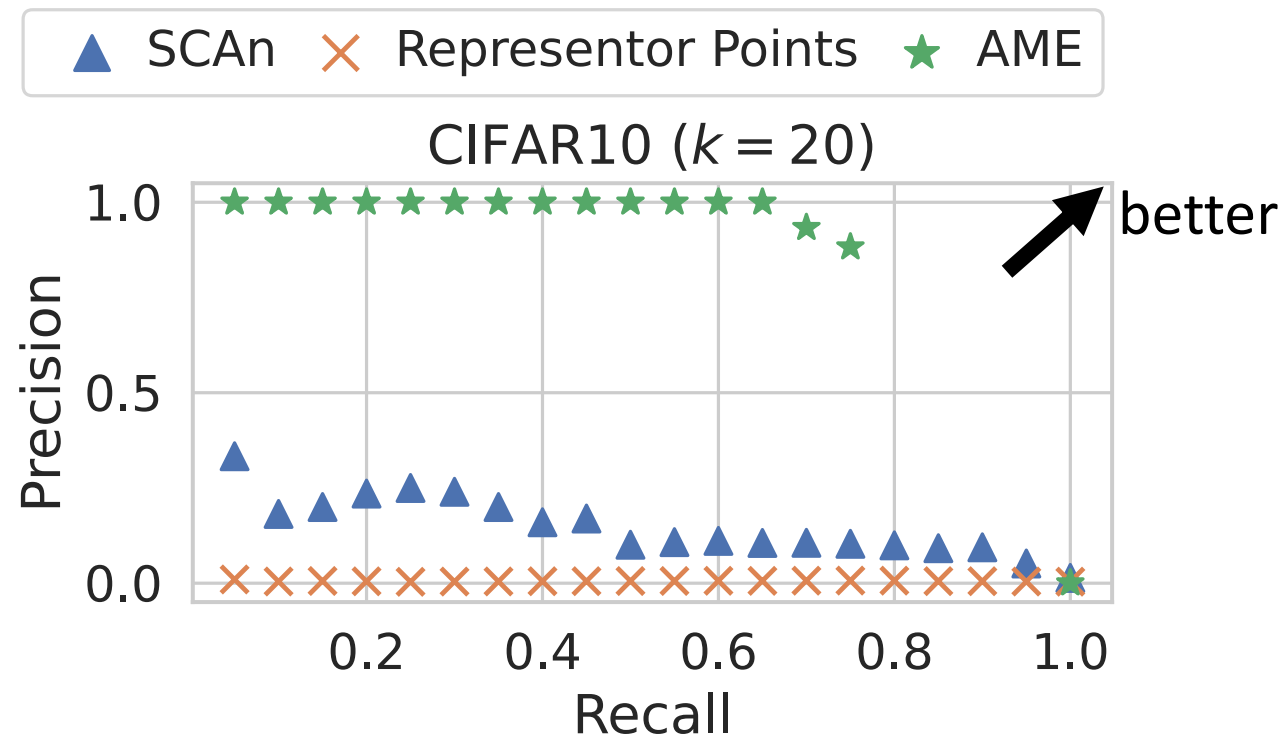
- Subsets are drawn by including each data point independently with probability $p \sim \mathcal{P}$
- $\mathcal{P} = \text{Uniform}(0,1)$, AME = SV, but incompatible with fast estimation
- Instead, $\mathcal{P} = \text{“discretized Uniform}(0,1)\text{”}$

Successfully Detects Poisoned Training Data

- We poisoned k training data points, estimate the AMEs, and select those with high AME

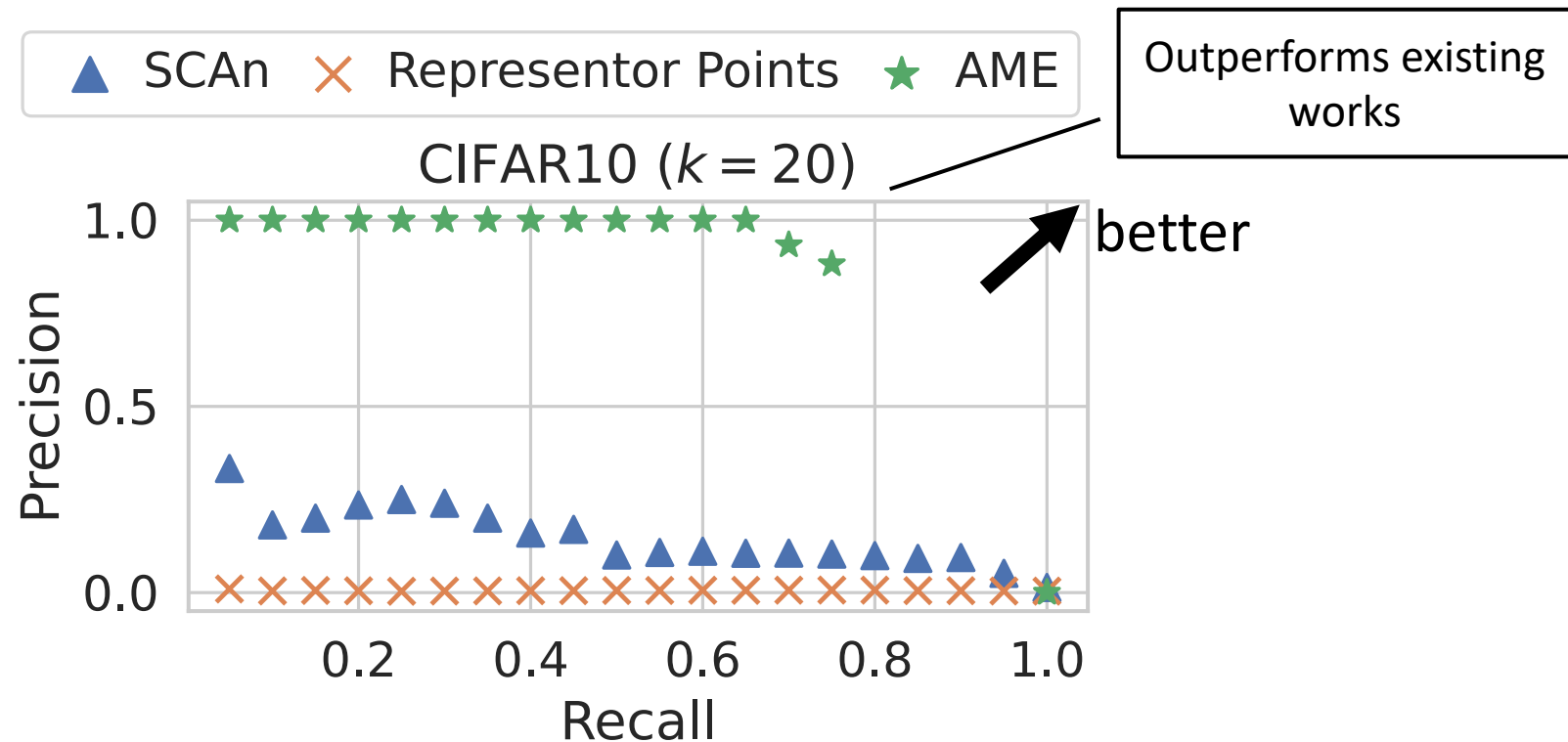
Successfully Detects Poisoned Training Data

- We poisoned k training data points, estimate the AMEs, and select those with high AME



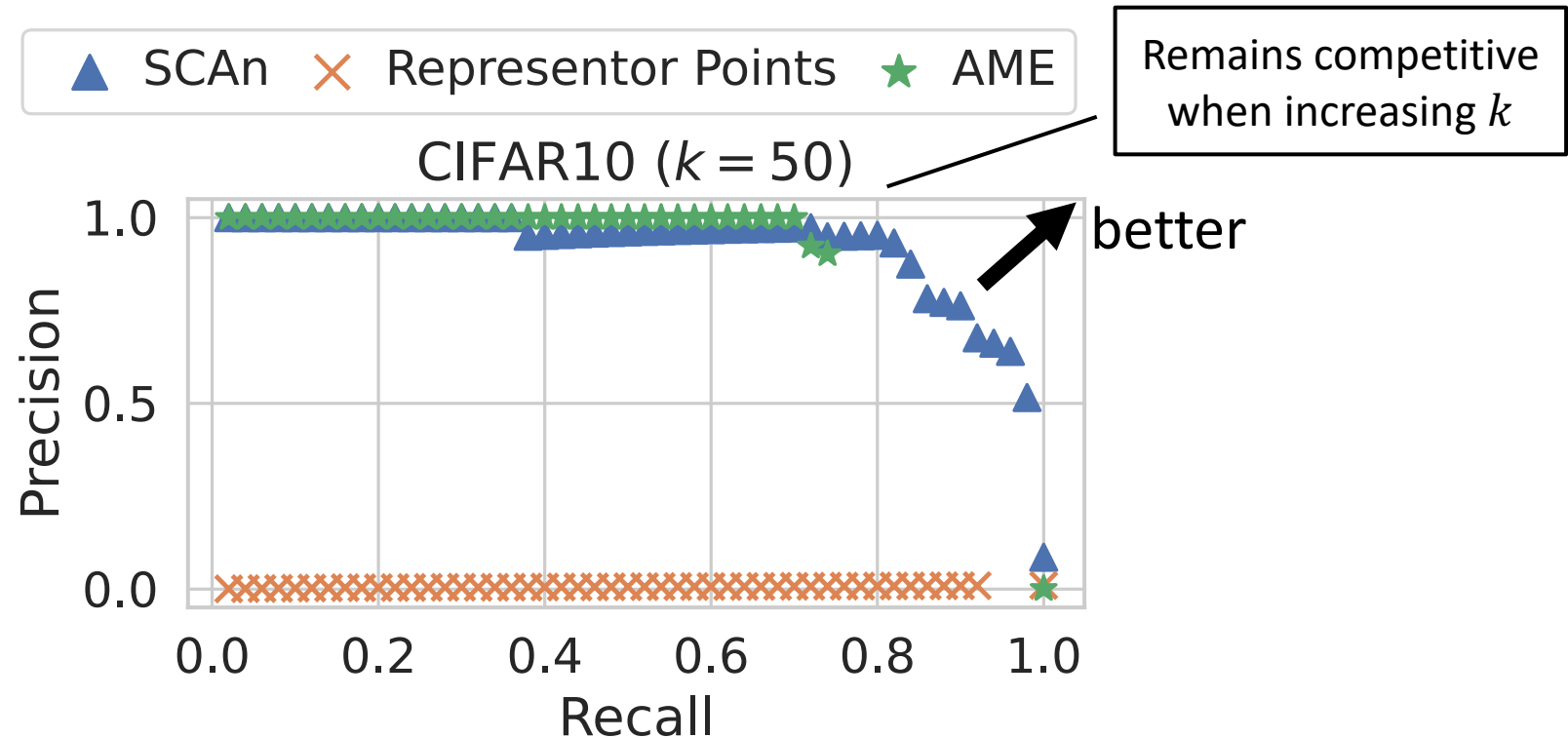
Successfully Detects Poisoned Training Data

- We poisoned k training data points, estimate the AMEs, and select those with high AME



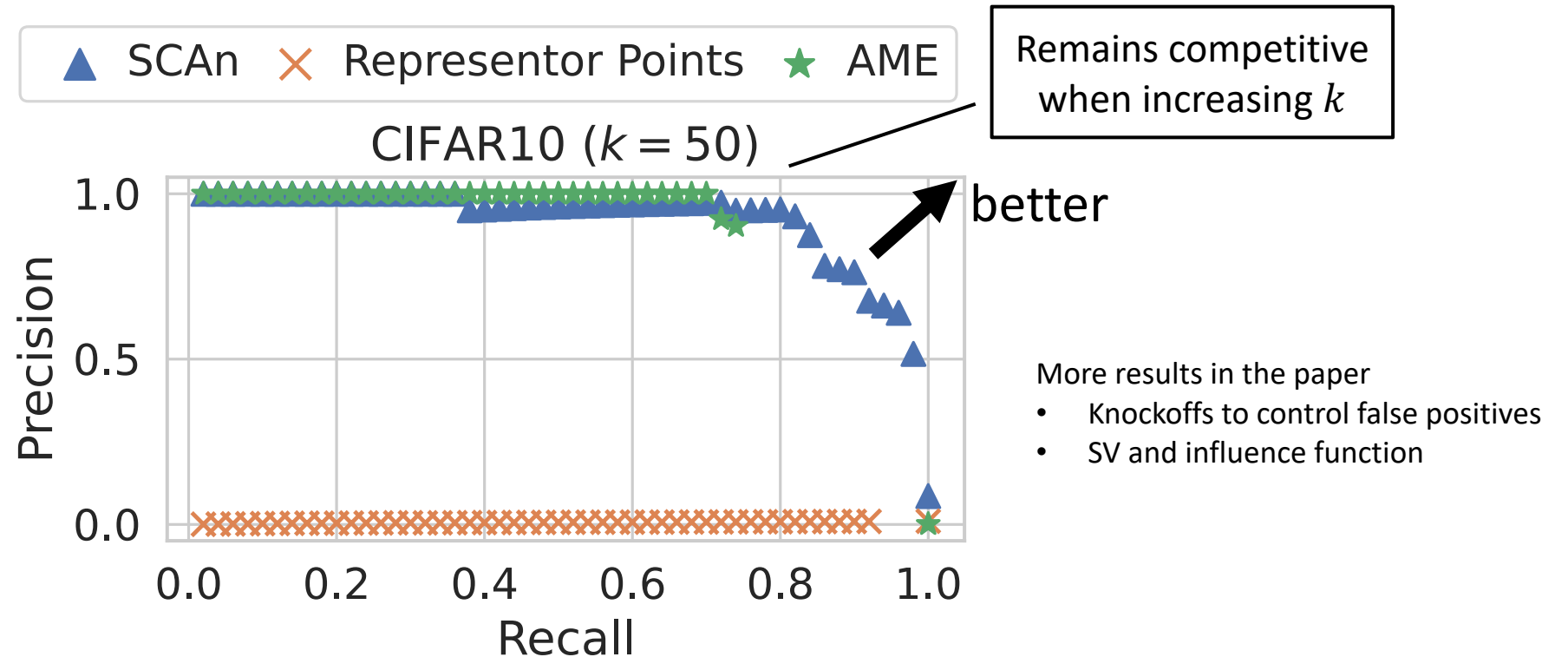
Result in poison detection

- We poisoned k training data points, estimate the AMEs, and select those with high AME

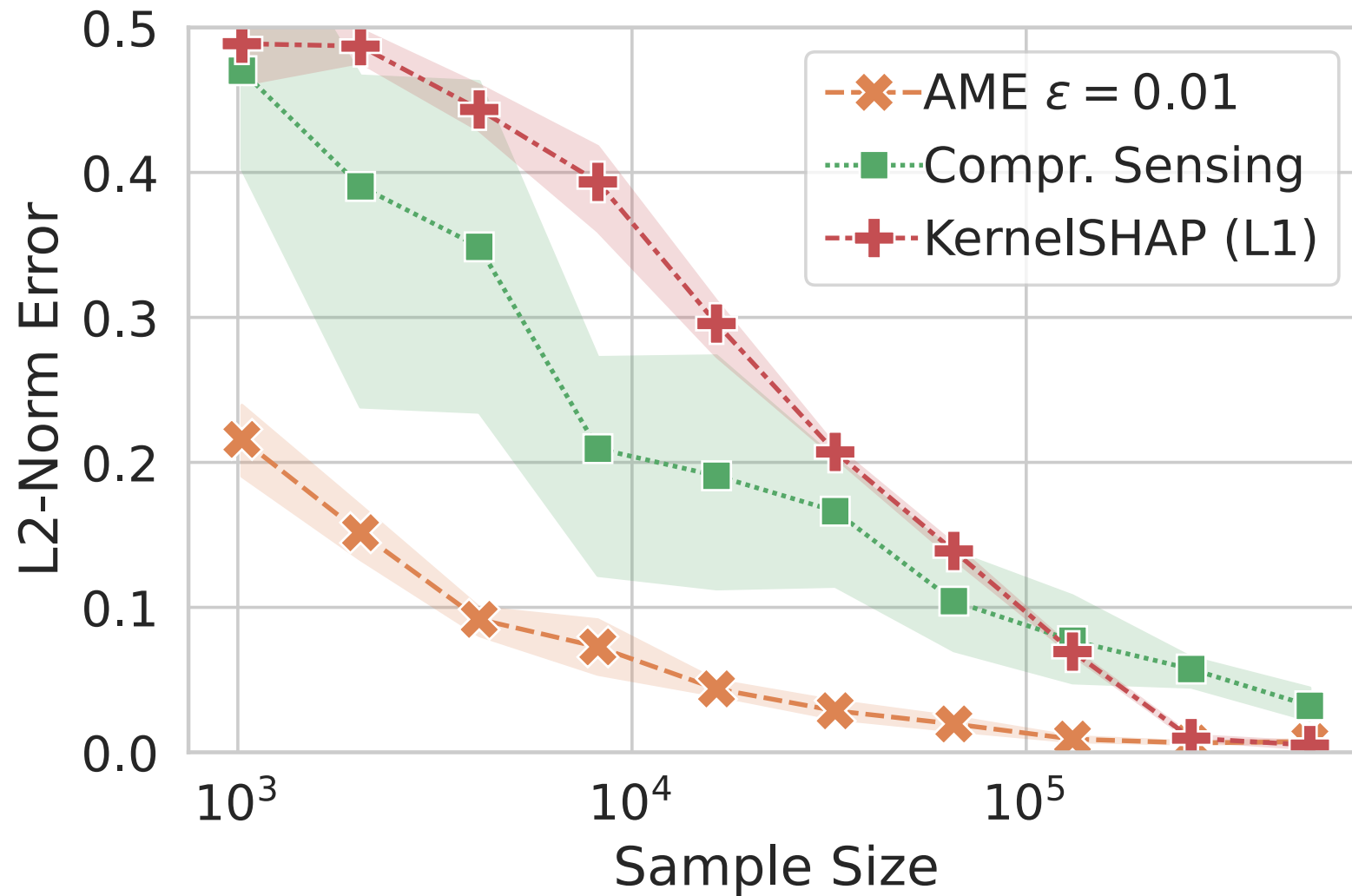


Result in poison detection

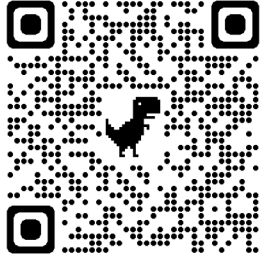
- We poisoned k training data points, estimate the AMEs, and select those with high AME



Faster estimator for SV

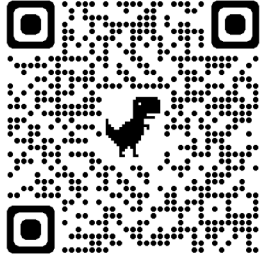


Conclusion



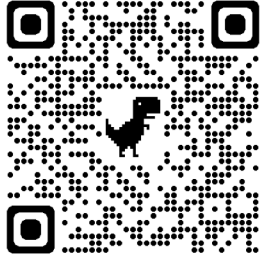
- A new method useful for model debugging and poison detection, assuming sparsity
 - Uses AME to measure effect of training data points
 - Uses randomized experiments and LASSO to estimate AME efficiently
 - Faster estimator for SV

Conclusion



- A new method useful for model debugging and poison detection, assuming sparsity
 - Uses AME to measure effect of training data points
 - Uses randomized experiments and LASSO to estimate AME efficiently
 - Faster estimator for SV
- See our paper for more details, results and extensions:
 - Hierarchical design: simultaneously identify sources (collections of data points) and individual data points with high AME

Conclusion



- A new method useful for model debugging and poison detection, assuming sparsity
 - Uses AME to measure effect of training data points
 - Uses randomized experiments and LASSO to estimate AME efficiently
 - Faster estimator for SV
- See our paper for more details, results and extensions:
 - Hierarchical design: simultaneously identify sources (collections of data points) and individual data points with high AME
- Contact: jinkun.lin@nyu.edu; Poster: Hall E #936

Thank you!