

Tracing Data Errors Using View-Conditioned Causality

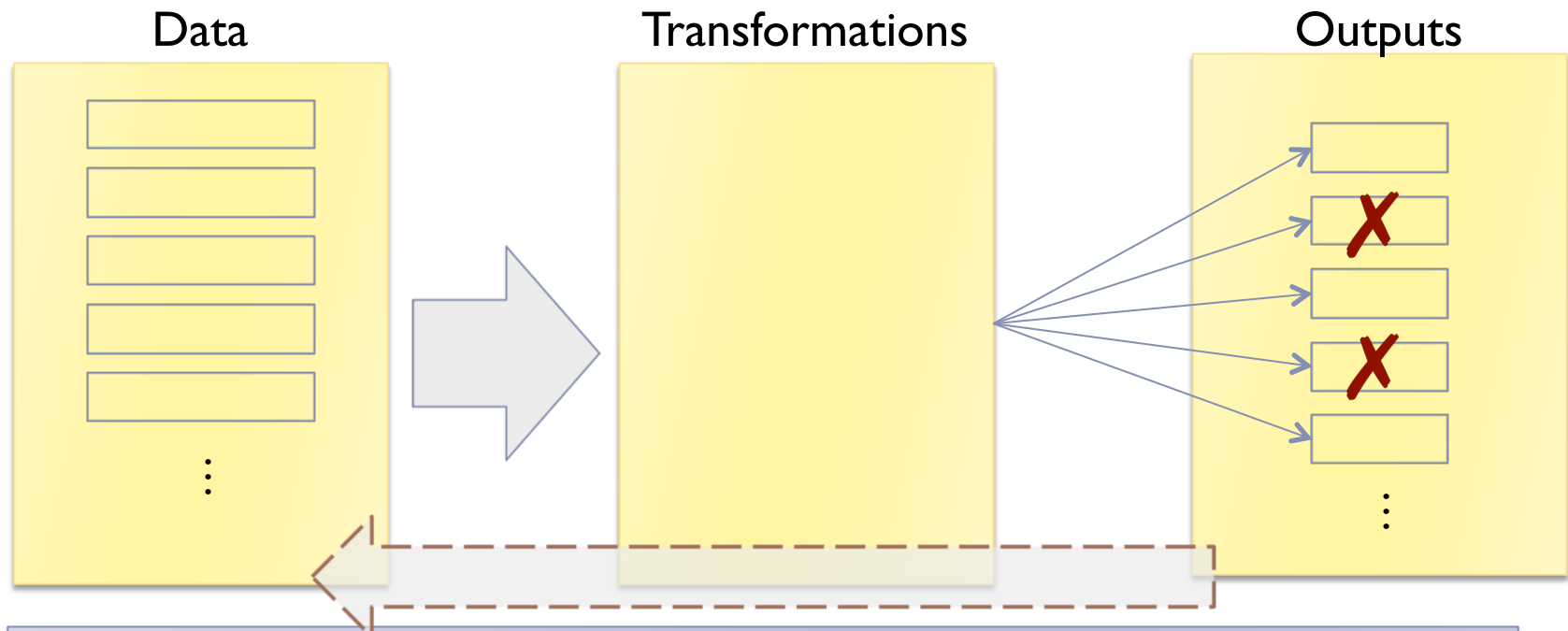
Alexandra Meliou^{*}

with Wolfgang Gatterbauer^{*}, Suman Nath[§], and Dan Suciu^{*}

^{}University of Washington, [§]Microsoft Research*



General Problem Setting

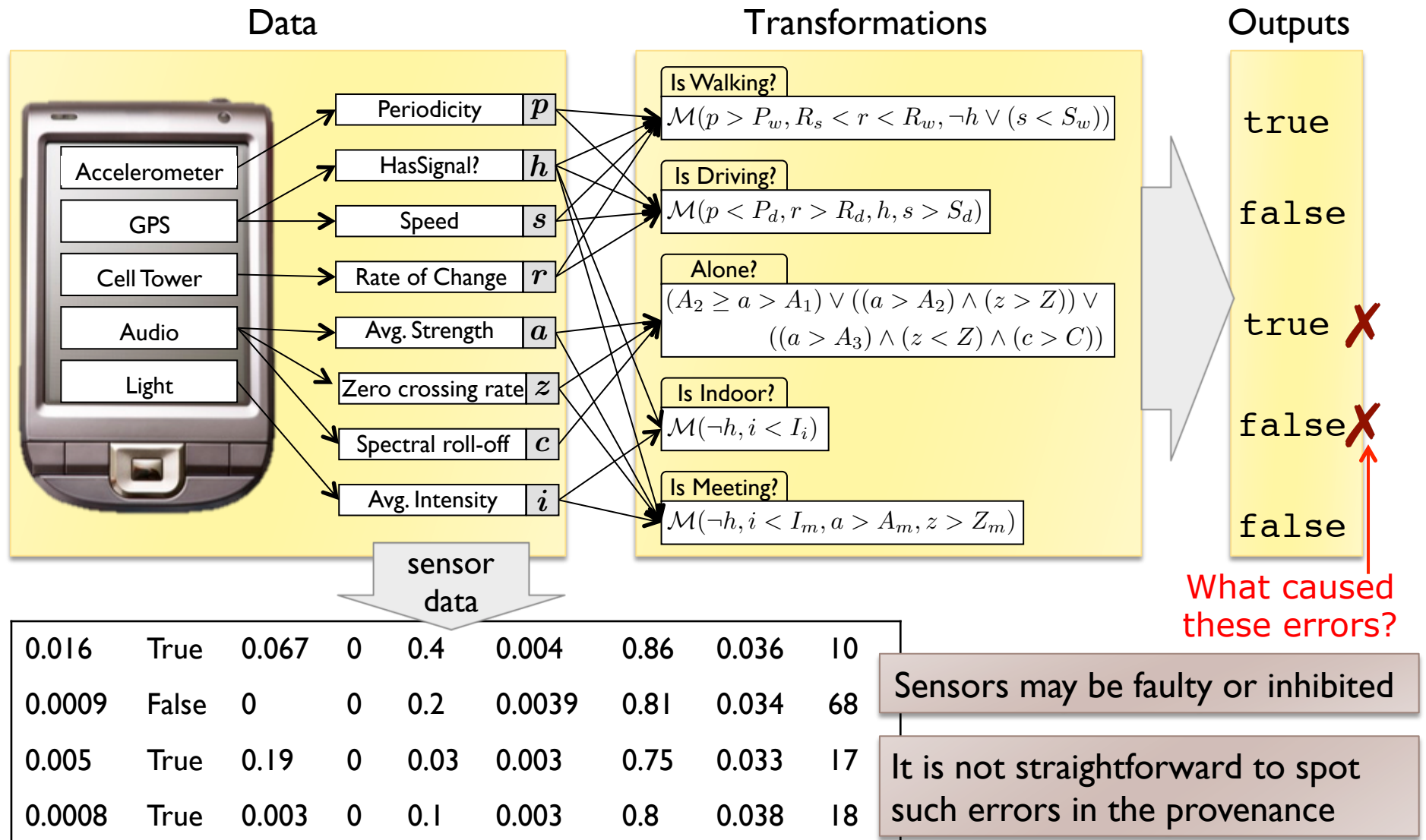


If one or more of the outputs are deemed erroneous, can we find the tuples in the base data responsible for that error?

Correcting those can fix even more potential errors in the output.

Provenance helps narrow down the candidate tuples in the input data. The challenge is to identify the input tuples that can best explain the observed errors.

Focus: Context Aware Recommendations

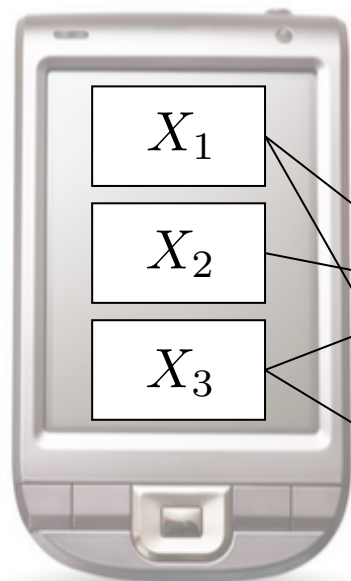


Contributions

- ▶ Introduce **view-conditioned causality** and **responsibility** for tracing errors in views and transformations to source data
 - ▶ The presence of errors is often obvious in the transformations but not the source data (*post-factum cleaning*)
- ▶ Non-trivial reduction of causality and responsibility computation to a **satisfiability** problem
- ▶ An optimized conversion algorithm that reduces the SAT problem size
- ▶ Illustration of effectiveness in a real-world classifier-based recommendation system using mobile sensor data
 - ▶ **High average precision, and almost 90% correction ratio in some cases**

Running Example

Input variables can be from a continuous or discrete domain



Φ_1

$$Z_1 = (X_1 < 5) \wedge (X_3 = 4) \vee \neg X_2$$

Φ_2

$$Z_2 = (X_1 > 2) \wedge (X_3 \geq 4)$$

Example:

Input

$$X_1 = 3$$

$$X_2 = \text{true}$$

$$X_3 = 4$$

Results in output

$$Z_1 = \text{true}$$

$$Z_2 = \underline{\text{true}}$$

error

But what if we know that the first classifier should evaluate to true, and the second to false?

Ground truth: $\hat{Z}_1 = \text{true}, \hat{Z}_2 = \underline{\text{false}}$

View-Conditioned Causality

Refer to the paper for the formal definitions

- ▶ A set of input variables is a **counterfactual cause**, if changing their values results in the correct output for all transformations, and the set is minimal.

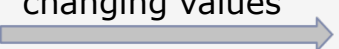
Example:

	Evaluate to:	Ground truth:	$X_1 = 3$
$Z_1 = (X_1 < 5) \wedge (X_3 = 4) \vee \neg X_2 = \text{true}$		$\hat{Z}_1 = \text{true}$	$X_2 = \text{true}$
$Z_2 = (X_1 > 2) \wedge (X_3 \geq 4)$	$= \text{true}$	$\hat{Z}_2 = \text{false}$	$X_3 = 4$

Counterfactual causes:

Change:

Gives output:

$\{X_1\}$ 

$X'_1 = 1$

$\{Z'_1, Z'_2\} = \{\hat{Z}_1, \hat{Z}_2\}$ **ground truth**

$\{X_2, X_3\}$ 

$\{X'_2, X'_3\} = \{\text{false}, 2\}$

$\{Z'_1, Z'_2\} = \{\hat{Z}_1, \hat{Z}_2\}$

View-Conditioned Causality

Refer to the paper for the formal definitions

- ▶ A variable is a **cause** if it is a part of a counterfactual cause
- ▶ If $X_i \cup \Gamma$ is a counterfactual cause, Γ is a **contingency** for X_i

Responsibility: $\rho_{X_i} = \frac{1}{1 + \min_{\Gamma} |\Gamma|}$ ← The smaller the contingency set, the higher the responsibility

Example:

	Evaluate to:	Ground truth:	$X_1 = 3$
$Z_1 = (X_1 < 5) \wedge (X_3 = 4) \vee \neg X_2 = \text{true}$		$\hat{Z}_1 = \text{true}$	$X_2 = \text{true}$
$Z_2 = (X_1 > 2) \wedge (X_3 \geq 4)$	$= \text{true}$	$\hat{Z}_2 = \text{false}$	$X_3 = 4$

Counterfactual causes:

$\{X_1\}$

$\{X_2, X_3\}$

Change:

$X'_1 = 1$

$\{X'_2, X'_3\} = \{\text{false}, 2\}$

Responsibility:

$\rho_{X_1} = 1$

$\rho_{X_2} = \rho_{X_3} = \frac{1}{2}$

causes

Our Approach to Post-Factum Cleaning

- ▶ Compute all causes and rank them by their responsibility.
 - ▶ Use the ranking as an indicator for error tracing
- ▶ But: Computing responsibility is hard for general Boolean formulas [Eiter et al. 2002], and even for conjunctive queries [PVLDB 2010]
- ▶ Transform causality into a satisfiability problem and use highly optimized SAT solvers, which are **very efficient in practice**
 - ▶ We explain how we do this in 4 main steps

Reduction to SAT

1. Map continuous input to Boolean partition variables

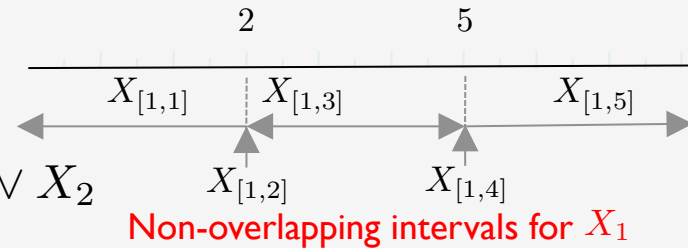
Example (cont.):

$$\Phi_1 : Z_1 = (X_1 < 5) \wedge (X_3 = 4) \vee \neg X_2$$

$$= (X_{[1,1]} \vee X_{[1,2]} \vee X_{[1,3]}) \wedge X_{[3,2]} \vee X_2$$

$$\Phi_2 : Z_2 = (X_1 > 2) \wedge (X_3 \geq 4)$$

$$= (X_{[1,3]} \vee X_{[1,4]} \vee X_{[1,5]}) \wedge (X_{[3,2]} \vee X_{[3,3]})$$



2. When the intervals are non-overlapping, we can easily model their correlation with a constraint

$$\Psi_i = \left(\bigvee_j X_{[i,j]} \right) \left(\bigwedge_{j < l} (\neg X_{[i,j]} \vee \neg X_{[i,l]}) \right) \quad \Psi = \bigwedge \Psi_i$$

At least one is true + No two are true together = exactly one is true

Example (cont.): $\Psi_3 = (X_{[3,1]} \vee X_{[3,2]} \vee X_{[3,3]})$

$$\wedge (\neg X_{[3,1]} \vee \neg X_{[3,2]}) \wedge (\neg X_{[3,1]} \vee \neg X_{[3,3]}) \wedge (\neg X_{[3,2]} \vee \neg X_{[3,3]})$$

Reduction to SAT

Running Example:

$$\begin{aligned}\Phi_1 : Z_1 &= (X_1 < 5) \wedge (X_3 = 4) \vee \neg X_2 \\ \Phi_2 : Z_2 &= (X_1 > 2) \wedge (X_3 \geq 4)\end{aligned}$$

Ground truth:

$\hat{Z}_1 = \text{true}$

$\hat{Z}_2 = \text{false}$

Input values:

$X_1 = 3$

$X_2 = \text{true}$

$X_3 = 4$

3. a. **Construct a Boolean formula whose satisfying assignments produce the correct output**

$$\hat{\Phi} = \left(\bigwedge_{i: \hat{Z}_i = \text{T}} \Phi_i \right) \wedge \left(\bigwedge_{i: \hat{Z}_i = \text{F}} \neg \Phi_i \right)$$

Example (cont.):

$$\hat{\Phi} = \Phi_1 \wedge \neg \Phi_2$$

All satisfying assignments of $\hat{\Phi}$ cause each Φ_i to evaluate to its ground truth

- b. **Construct a Boolean formula whose satisfying assignments satisfy $\hat{\Phi}$, and also change the value of X_i**

$$\Phi_{\text{SAT}} = \neg \hat{\Phi} [\theta(\mathbf{X}_{[i]})] \wedge \hat{\Phi} [\neg \theta(\mathbf{X}_{[i,j]})] \wedge \Psi [\neg \theta(\mathbf{X}_{[i,j]})] \quad (\text{hard constraint})$$

Example (cont.): X_1 is a cause iff the following formula is satisfiable:

$$\Phi_{\text{SAT}} = \neg \hat{\Phi} [\{X_{[1,1]}, X_{[1,2]}, X_{[1,3]}, X_{[1,4]}, X_{[1,5]}\} = \{\text{F}, \text{F}, \text{T}, \text{F}, \text{F}\}]$$

Current assignment of X_1

$$\wedge \hat{\Phi} [X_{[1,3]} = \text{F}] \wedge \Psi [X_{[1,3]} = \text{F}]$$

Negate current assignment of X_1

Computing Responsibility with MaxSAT

Running Example:

$$\begin{aligned}\Phi_1 : Z_1 &= (X_1 < 5) \wedge (X_3 = 4) \vee \neg X_2 \\ \Phi_2 : Z_2 &= (X_1 > 2) \wedge (X_3 \geq 4)\end{aligned}$$

Ground truth:

$$\hat{Z}_1 = \text{true}$$

$$\hat{Z}_2 = \text{false}$$

Input values:

$$X_1 = 3$$

$$X_2 = \text{true}$$

$$X_3 = 4$$

4. Construct “soft” constraints to find minimum contingency set

$$\Phi_\theta = \bigwedge_{\theta(X_{[i,j]})=\text{T}} X_{[i,j]} \bigwedge_{\theta(X_{[i,j]})=\text{F}} \neg X_{[i,j]} \quad (\text{soft constraint})$$

Example (cont.):

$$\begin{aligned}\Phi_\theta = & \neg X_{[1,1]} \wedge \neg X_{[1,2]} \wedge X_{[1,3]} \wedge \neg X_{[1,4]} \wedge \neg X_{[1,5]} \\ & \wedge X_2 \wedge \neg X_{[3,1]} \wedge \neg X_{[3,2]} \wedge X_{[3,3]}\end{aligned}$$

A partial MaxSAT solver tries to satisfy as many conjuncts of the soft constraint as possible, and thus produces an assignment as similar to the given one as possible

Minimum contingency

Experimental Setup

- ▶ Three individuals using our context-aware recommendation system on their mobile devices over a period of 3 weeks
- ▶ Dataset:
 - ▶ 800 different instances of user activity
 - ▶ 150 total hours of data during the 3 weeks
- ▶ The users recorded erroneous outputs, as well as whether sensors happened to be inhibited
- ▶ SAT reduction implemented in Java, output exported in standard DIMACS CNF and WCFN formats
- ▶ MiniSat (<http://minisat.se/>) and MiniMaxSat ([Heras et al. 2008]) solvers

Average Precision

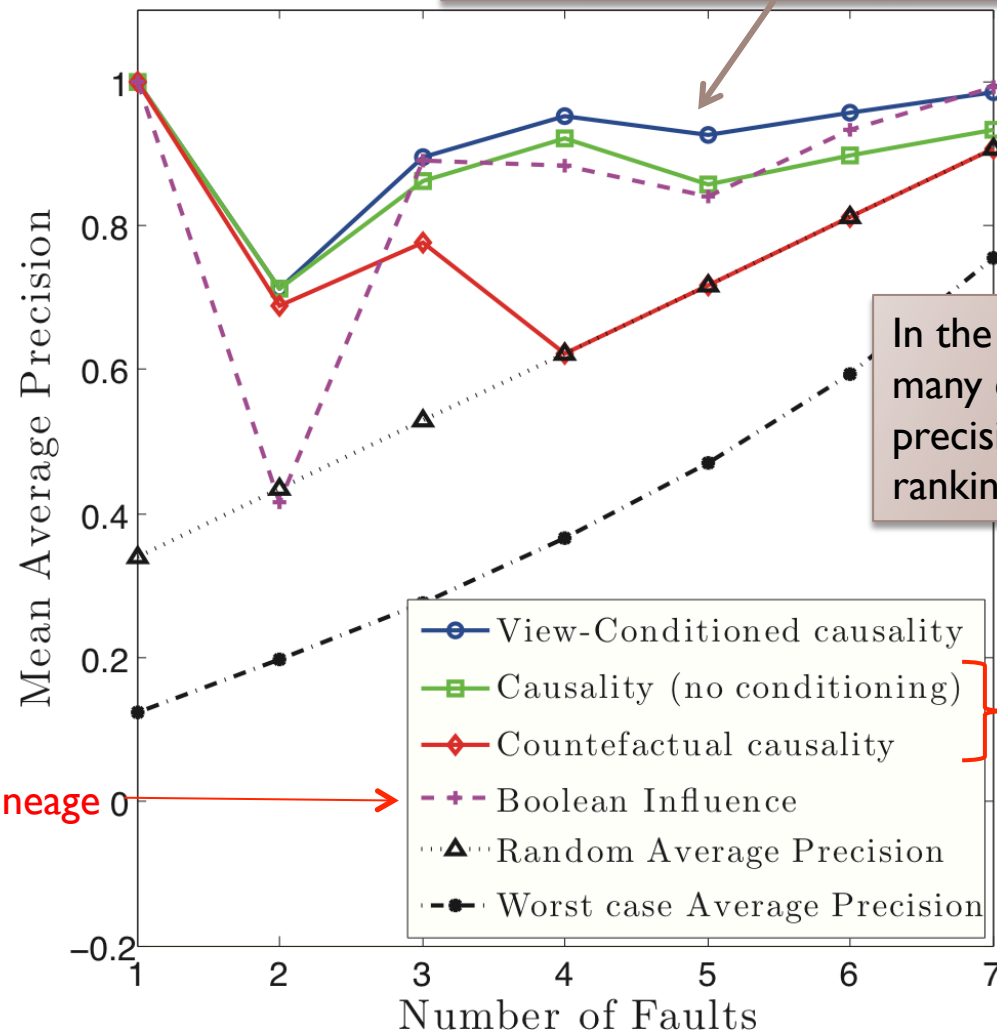
800 different instances
5 sensory inputs
8 extracted features (variables)
3 users

Average precision is a metric of quality of a ranking.

If all erroneous variables are ranked first, then average precision is 1.

Static analysis of lineage →

View-Conditioned causality produces more accurate error rankings than other approaches



In the presence of many errors the avg precision of all rankings increases

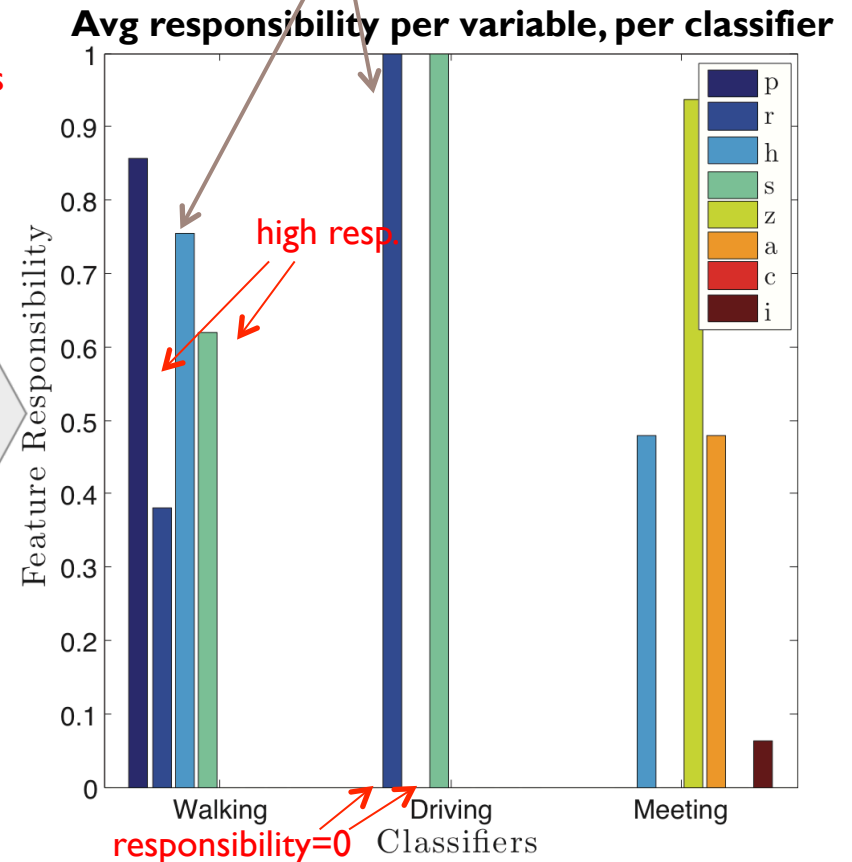
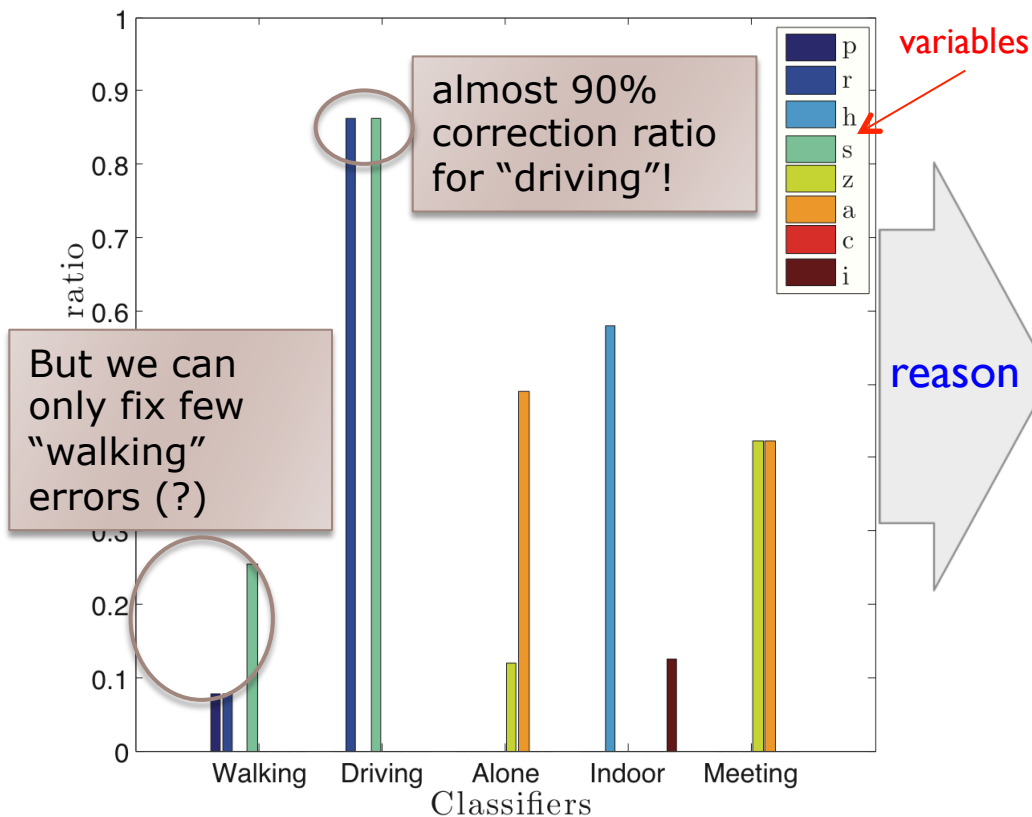
Simpler causality schemes

Corrections

We select the highest responsibility variable, remove it from the evaluation of all classifiers, and record the portion of errors that get corrected per classifier

Driving has reliable features (low responsibility), means they are almost never causes of error

Walking has no reliable features

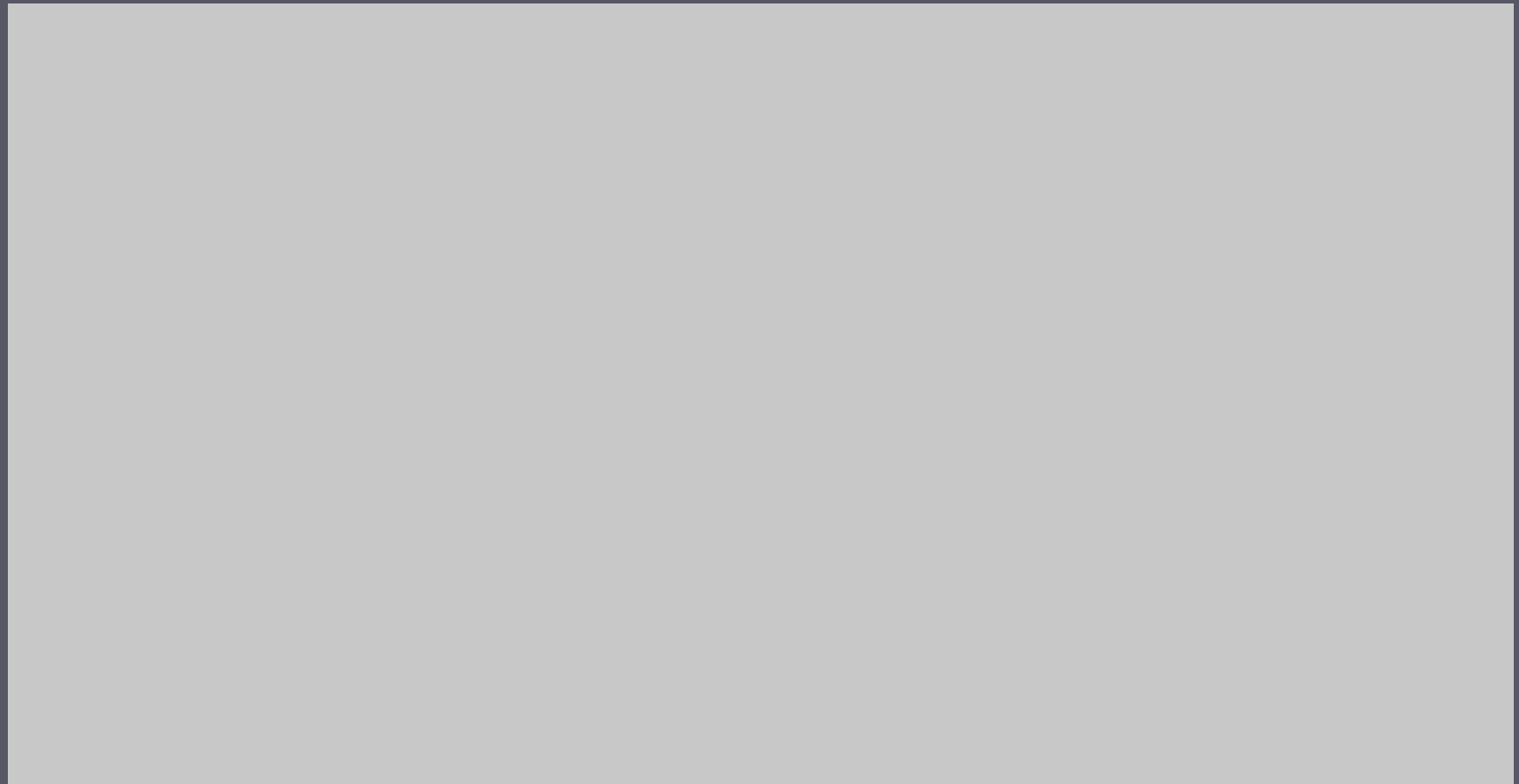


Conclusions

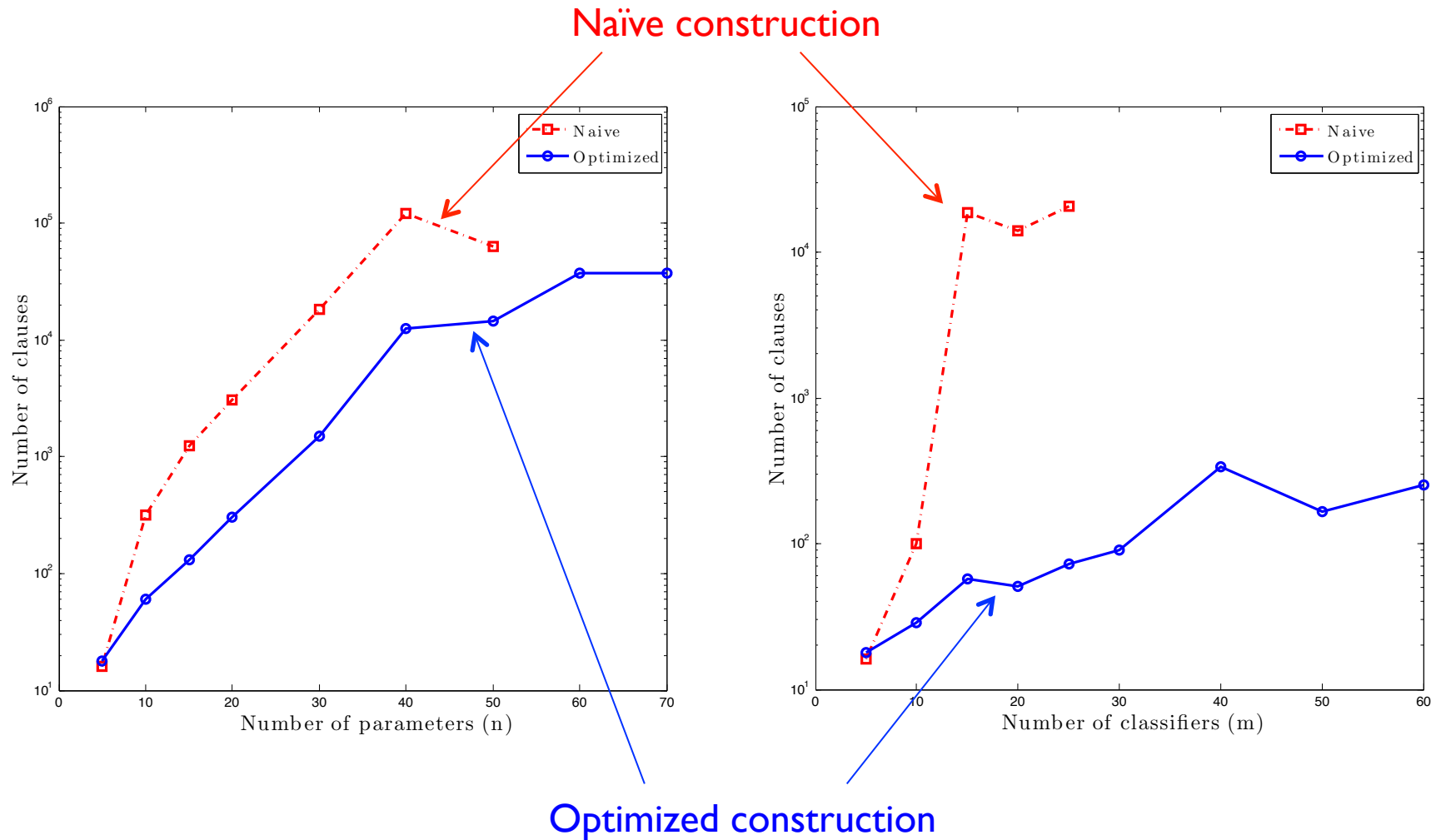
- ▶ Defined view-conditioned causality (VCC) and demonstrated its effectiveness in post-factum cleaning
 - ▶ Results show that VCC successfully identifies causes of error
- ▶ Described a non-trivial reduction to a satisfiability problem
- ▶ Also in the paper
 - ▶ Optimization of formula size (we achieve orders of magnitude improvement)
 - ▶ Scalability experiments
- ▶ Questions?



Additional Graphs



Improving the CNF Size



SAT Solver Runtime

