# Data Glitches = Constraint Violations – Empirical Explanations

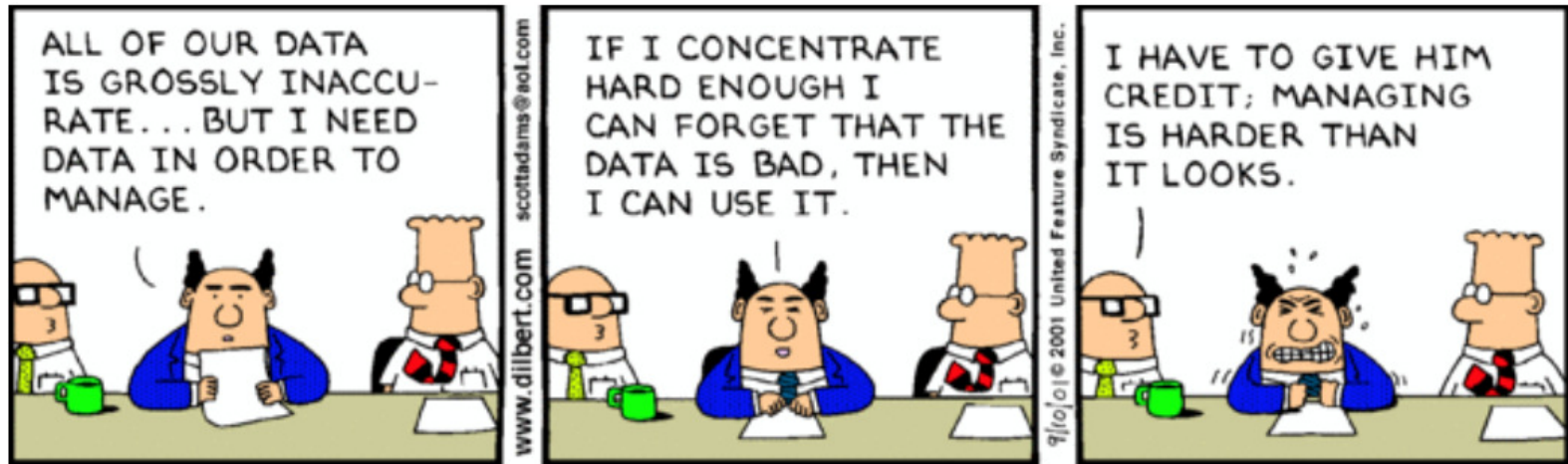**Divesh Srivastava**

**AT&T Labs-Research**

# What is a Glitch?



♦ A spaceman's word for irritating disturbances [Time, 23 Jul 1965].

   – "Something's gone wrong and you can't figure out what it is" [Daly].

# What is a Data Glitch?



♦ A data scientist's phrase for irritating data quality problems.
   – Data that has gone wrong and **can't be used as desired**.
   – Unusual data that **does not conform** to data quality expectations.

# What is an Integrity Constraint Violation?

♦ Integrity constraint: formal specification that data must satisfy.
  – **Semantic** (SSN unique for person) vs **syntactic** (NNN-NN-NNNN).
  – **Logical** (FD on 52wk low-high) vs **statistical** (# files within $3\sigma$ of $\mu$).

♦ Violation: data that does not satisfy specified integrity constraint.

# "Small Data" Quality: How Was It Achieved?

◆ Specify all domain knowledge as **integrity constraints** on data.

– **Reject updates** that do not preserve integrity constraints.

– Works well when the domain is very well understood and static.

# Data Quality: Impact of Big Data



♦ **Variety, variability** of data: one size does not fit all.

# Big Data



♦ Big data is different things to different people.
  – Volume, velocity, variety, variability, value, **veracity.**

# Big Data Quality: A Different Approach?

♦ Big data: integrity constraints cannot be always specified a priori.

- Data **variety** → complete domain knowledge is infeasible.
- Data **variability** → domain knowledge becomes obsolete.
- Too much rejected data → "small" data. ☺

# Big Data Quality: A Different Approach?

♦ Big data: integrity constraints cannot be always specified a priori.
  – Data **variety** → complete domain knowledge is infeasible.
  – Data **variability** → domain knowledge becomes obsolete.

♦ Solution: let the data speak for itself.
  – Learn (simple) **integrity constraints / models** from the data.
  – Identify **violations** of the learned constraints.
  – Learn (complex) **empirical explanations** of the identified violations.
  – Declare **glitches** = constraint violations – empirical explanations.

# In This Talk

◆ Big data: integrity constraints cannot be always specified a priori.

  – Data **variety** → complete domain knowledge is infeasible.

  – Data **variability** → domain knowledge becomes obsolete.

◆ Solution: let the data speak for itself.

  – Learn (simple) **integrity constraints / models** from the data.

  – Identify **violations** of the learned constraints.

  – Learn (complex) **empirical explanations** of the identified violations.

  – Declare **glitches** = constraint violations – empirical explanations.

# Outline

♦ Introduction.

♦ What is an empirical explanation?

♦ Unsupervised learning of empirical explanations.

# What is an Empirical Explanation?

| ID | Status | Phone | Dept. | Rm. | Super_ID |
|---|---|---|---|---|---|
| **ID_5** | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | **New Hire** | 1AAA3608776 | D2300 | D284 | **ID_5** |
| ID_8 | **New Hire** | 1AAA3608776 | D2300 | B106 | **ID_5** |

♦ Data does not conform to expectation of "phone # uniqueness".

– Explanation = "new hires can have same phone # as supervisor".

– Explanation **can be learned from the data** → empirical explanation.

# What is an Empirical Explanation? ✓

| ID | Status | Phone | Dept. | Rm. | Super_ID |
|---|---|---|---|---|---|
| ID_10 | Active | 1AAA3605519 | D8000 | **A132** | ID_13 |
| ID_11 | Active | 1AAA3605519 | D8000 | **A132** | ID_13 |
| ID_12 | Active | 1AAA3605519 | D8000 | **A132** | ID_13 |

◆ Data does not conform to expectation of "phone # uniqueness".

- Explanation = "employees in same room can have same phone #".
- Is this an empirical explanation?

# What is an Empirical Explanation?

| ID | Status | Phone | Dept. | Rm. | Super_ID |
|------|-----------|-------------|--------|---------|----------|
| ID_1 | Active | 1AAA3600000 | D4000 | ------- | ID_4 |
| ID_2 | --------- | 1AAA3600000 | ------- | ------- | ------ |
| ID_3 | Active | 1AAA3600000 | D2200 | E260 | ID_6 |

◆ Data does not conform to expectation of "phone # uniqueness".

– No empirical explanation is discernible.

– May be a **data glitch**. ☺

# What is an Empirical Explanation?



- ◆ Data does not conform to expectation of "FD on 52wk low-high".
  - – Explanation = "52 wk low-high definitions differ between sources".
  - – Is this an empirical explanation?

# What is an Empirical Explanation?



◆ Data does not conform to (statistical) expectation of "≤ 3σ of μ".

  – No empirical explanation is discernible; could it be a data glitch?

# What is an Empirical Explanation?



- ◆ Data does not conform to (statistical) expectation of "≤ 3σ of μ".
  - – Empirical explanation = "Fewer taxi trips during high wind speeds".
  - – An empirical explanation may **involve multiple data sets**.

# Outline

♦ Introduction.

♦ What is an empirical explanation?

♦ Unsupervised learning of empirical explanations.
  – Using spatio-temporal topological features [CD+16].
  – Using statistical signatures [DLS14].

# Using Spatio-Temporal Features: Problem



◆ Problem: Find data sets with **correlated spatio-temporal outliers**.

# Using Spatio-Temporal Features: Alternatives



◆ Traditional approaches: Pearson's correlation, DTW, etc.

– Miss **relationships that occur only at certain times / locations**, e.g., most of the time, # of taxi trips and wind speed are not related.

# Using Spatio-Temporal Features: Challenges



◆ Finding correlated spatio-temporal outliers is challenging.

– Big data sets, at different spatio-temporal resolutions.

– Combinatorial # of possible correlations to evaluate.

# Using Spatio-Temporal Features: Solution



♦ Solution: the **Data Polygamy** framework [CD+16].

– Constraint violations = topological features (e.g., peaks, valleys).

– Empirical explanations = significant (not a coincidence) correlations.

# Interesting Relationships Discovered

♦ Data sets: NYC urban, NYC open data.

♦ Weather and vehicle collisions.

    – Strong correlation between heavy rainfall and motorist fatalities.

    – No significant relationship between rainfall and vehicle collisions.

♦ Weather and taxi availability.

    – Strong correlation between heavy rainfall and number of taxis.



**Intelligencer**

## Why You Can't Get a Taxi When It's Raining

By **Annie Lowrey**    Follow @AnnieLowrey

Good luck, lady. Photo: Jacobs Stock Photography/Getty Images

It's pouring rain. You're running late. You desperately want to take a cab to the office. But, of course, there are none to be found. Happens all the time, right? Right, says science — or, to be specific, a new and exhaustive economic analysis of New York City taxi rides and Central Park meteorological data.

# Outline
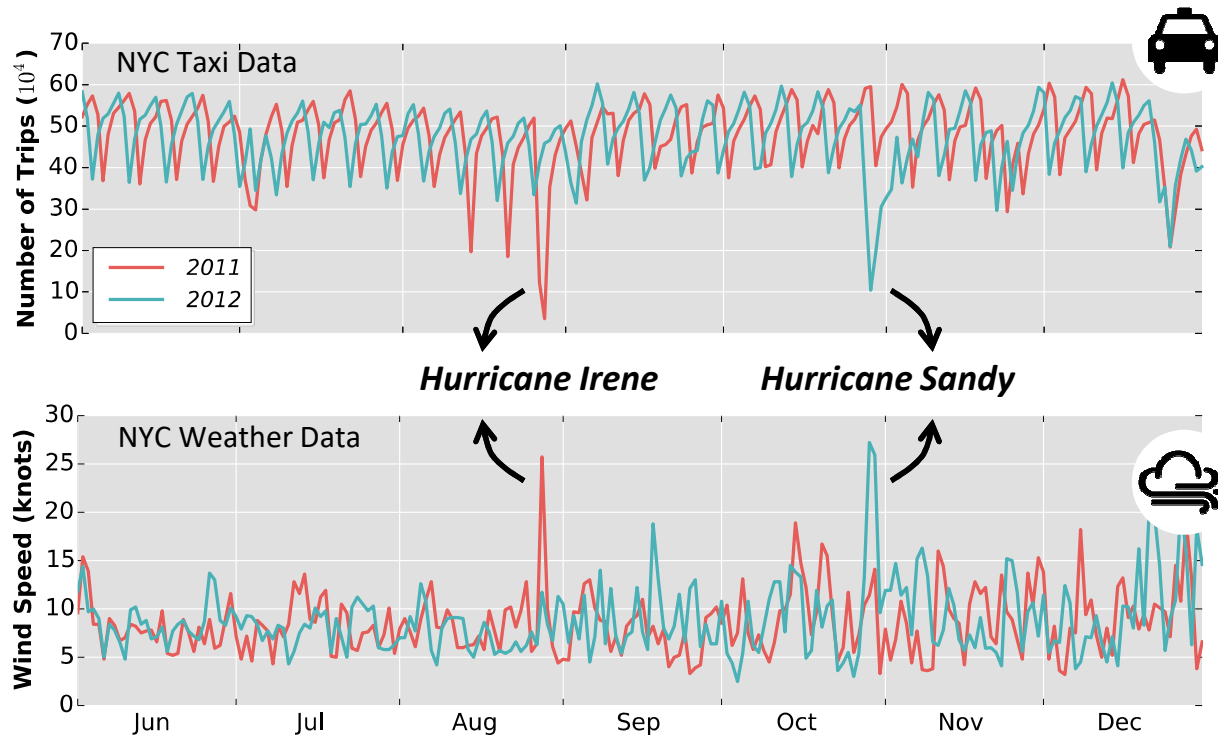
♦ Introduction.

♦ What is an empirical explanation?

♦ Unsupervised learning of empirical explanations.

   – Using spatio-temporal topological features [CD+16].

   – Using statistical signatures [DLS14].
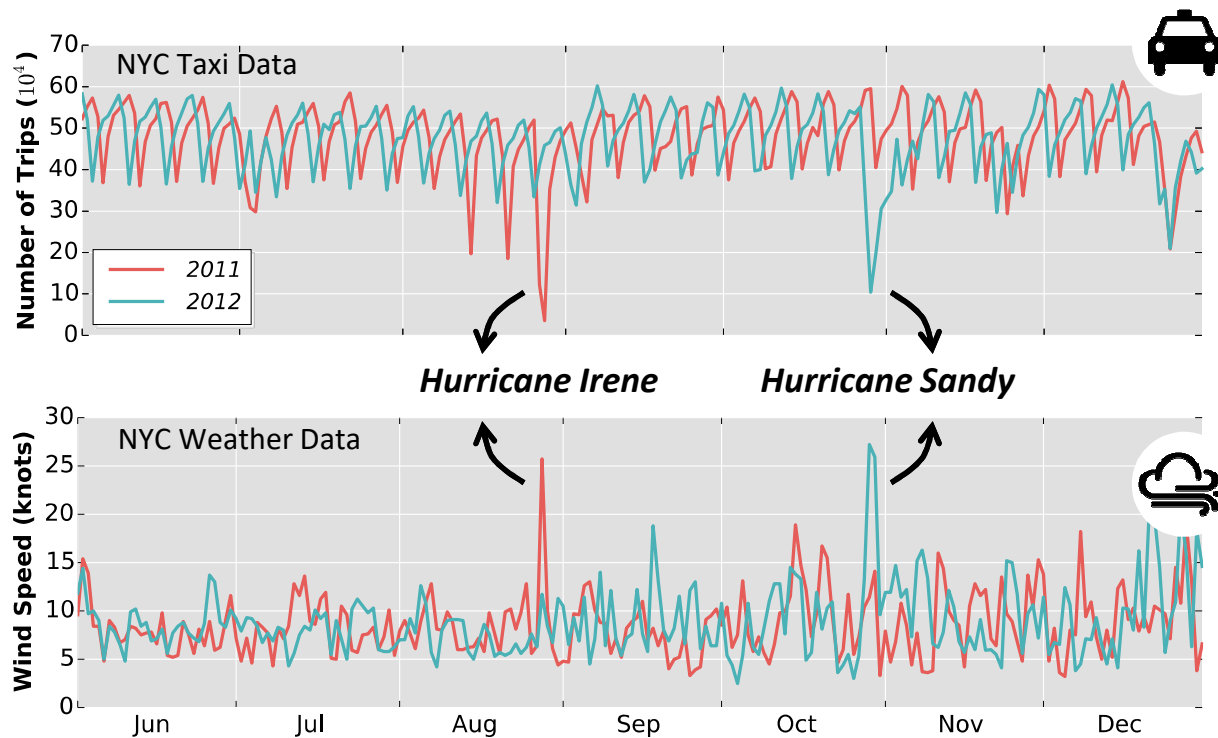
# Using Statistical Signatures: Problem

| ID | Status | Phone | Dept | Rm. | Super |
|---|---|---|---|---|---|
| **ID_5** | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | **New Hire** | 1AAA3608776 | D2300 | D284 | **ID_5** |
| ID_8 | **New Hire** | 1AAA3608776 | D2300 | B106 | **ID_5** |

♦ Problem: Find **statistically significant explanations** of violations.

  – Needed because of incomplete, obsolete domain knowledge.

# Using Statistical Signatures: Overview

| Apply data quality constraint | Identify **suspicious set** A |
|---|---|

| Compute **propensity signatures** of values *v* in A | Establish **statistical significance** of each *v* in A using **crossover subsampling** | Automatically generate **empirical explanations** using significant *v* |
|---|---|---|

Domain Expert

| Repair or Release | Refine domain knowledge |
|---|---|

26

# Using Statistical Signatures: Step 1

Data D

Suspicious set A: duplicate phone number

A'

Good Data

A

Suspicious Data

| ID | Status | Phone | Dept | Rm. | Super |
|------|----------|-------------|-------|------|-------|
| ID_5 | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | New Hire | 1AAA3608776 | D2300 | D284 | ID_5 |
| ID_8 | New Hire | 1AAA3608776 | D2300 | B106 | ID_5 |

◆ Apply constraint on D, identify violations (suspicious set) A.

◆ For each value v in A, compute **propensity signatures** in A and A'.

  – $s_A$(New Hire) = {0.67, 0.0, 0.0, 0.0, 0.0, 0.0}

  – $s_{A'}$(New Hire) = {0.05, 0.0, 0.0, 0.0, 0.0, 0.0}

# Using Statistical Signatures: Step 1

Data D

Suspicious set A: duplicate phone number



| ID | Status | Phone | Dept | Rm. | Super |
|---|---|---|---|---|---|
| **ID_5** | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | **New Hire** | 1AAA3608776 | D2300 | D284 | **ID_5** |
| ID_8 | **New Hire** | 1AAA3608776 | D2300 | B106 | **ID_5** |

Good Data — A'

Suspicious Data — A

◆ Apply constraint on D, identify violations (suspicious set) A.

◆ For each value v in A, compute **propensity signatures** in A and A'.

   – $s_A(\text{ID\_5}) = \{0.33, 0.0, 0.0, 0.0, 0.0, 0.67\}$

   – $s_{A'}(\text{ID\_5}) = \{0.02, 0.0, 0.0, 0.0, 0.0, 0.05\}$

# Using Statistical Signatures: Step 1

Data D

Suspicious set A: duplicate phone number

| ID | Status | Phone | Dept | Rm. | Super |
|------|----------|-------------|-------|------|-------|
| **ID_5** | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | **New Hire** | 1AAA3608776 | D2300 | D284 | **ID_5** |
| ID_8 | **New Hire** | 1AAA3608776 | D2300 | B106 | **ID_5** |

A'
Good Data

A
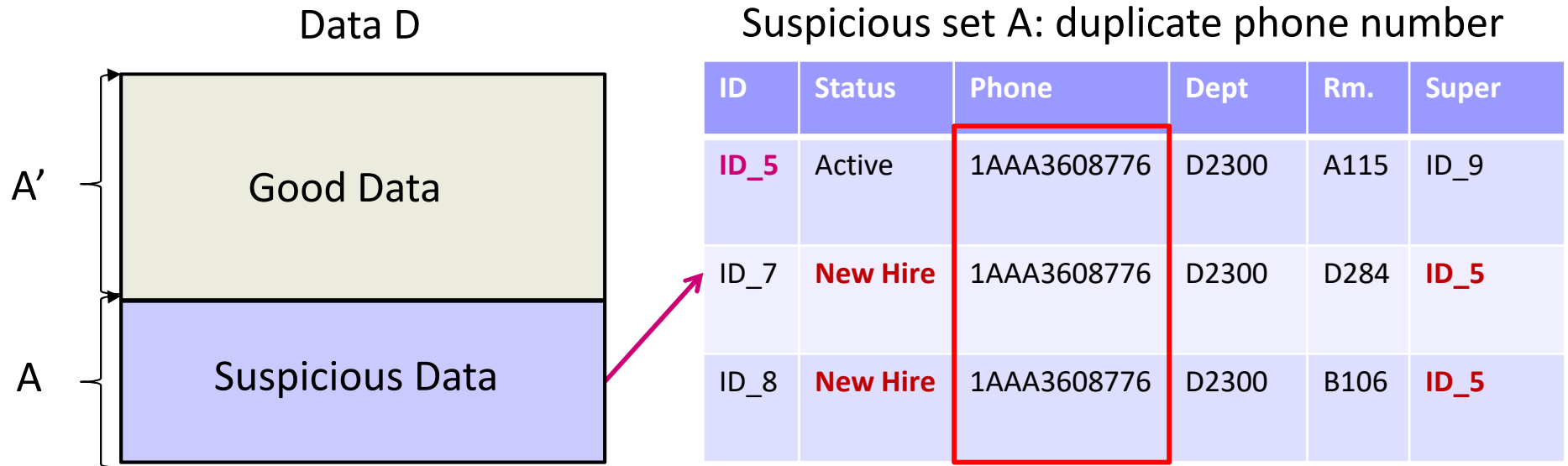Suspicious Data
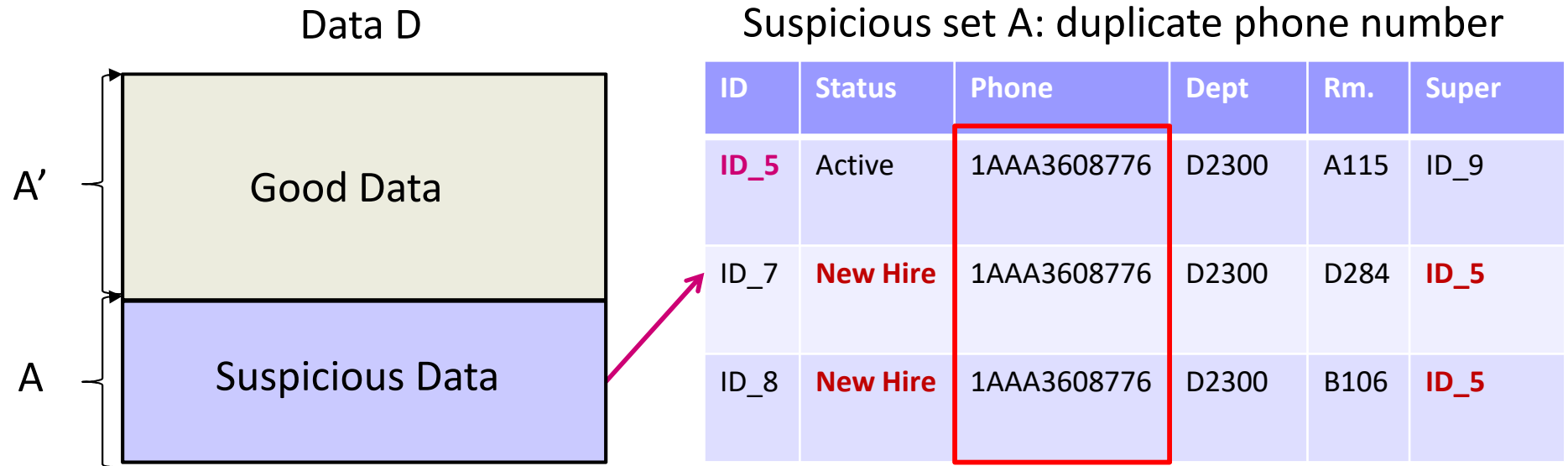
◆ Apply constraint on D, identify violations (suspicious set) A.

◆ For each value v in A, compute **propensity signatures** in A and A'.
  – Does value v have a **"sufficiently different" signature** in A vs A'?

29

# Using Statistical Signatures: Step 1

Data D

Suspicious set A: duplicate phone number



| ID | Status | Phone | Dept | Rm. | Super |
|------|----------|-------------|-------|------|-------|
| ID_5 | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | New Hire | 1AAA3608776 | D2300 | D284 | ID_5 |
| ID_8 | New Hire | 1AAA3608776 | D2300 | B106 | ID_5 |

◆ Apply constraint on D, identify violations (suspicious set) A.

◆ For each value v in A, compute **propensity signatures** in A and A'.

  − $s_A$(New Hire) = {0.67, 0.0, 0.0, 0.0, 0.0, 0.0}

  − $s_{A'}$(New Hire) = {0.05, 0.0, 0.0, 0.0, 0.0, 0.0}

# Using Statistical Signatures: Step 2

<table>
<tr><td colspan="6" style="text-align:center">Good Data</td></tr>
<tr><th>ID</th><th>Status</th><th>Phone</th><th>Dept</th><th>Rm.</th><th>Super</th></tr>
<tr><td>ID_5</td><td>Active</td><td>1AAA3608776</td><td>D2300</td><td>A115</td><td>ID_9</td></tr>
<tr><td>ID_7</td><td>New Hire</td><td>1AAA3608776</td><td>D2300</td><td>D284</td><td>ID_5</td></tr>
<tr><td>ID_8</td><td>New Hire</td><td>1AAA3608776</td><td>D2300</td><td>B106</td><td>ID_5</td></tr>
</table>

A'

A

Crossover Subsample

**Blocks from A'**

**Block from A**

- ◆ Goal: informative values that distinguish A from A'.
  - – Establish statistical significance using **crossover subsampling**.
  - – For an A block, sample A' blocks R times to create distribution.

# Using Statistical Signatures: Step 3

| ID | Status | Phone | Dept | Rm. | Super |
|---|---|---|---|---|---|
| **ID_5** | Active | 1AAA3608776 | D2300 | A115 | ID_9 |
| ID_7 | **New Hire** | 1AAA3608776 | D2300 | D284 | **ID_5** |
| ID_8 | **New Hire** | 1AAA3608776 | D2300 | B106 | **ID_5** |

♦ **Empirical explanation**: collection of all informative values for A.

– Learned in an **unsupervised manner**, e.g., {ID_5, New Hire}.

– Experts check empirical explanations, and decide on actions taken.

# Summary

- Big data quality: let the data speak for itself.
    - Learn simple constraints from the data sets, identify violations.
    - Learn complex **empirical explanations** within and across data sets.
    - Data glitches = constraint violations − empirical explanations.

- Benefits: statistically robust, computationally efficient cleaning.
    - Reduces statistical distortion due to unnecessary cleaning.
    - Addresses challenges due to **variety, variability** in big data.

- Just the beginning, a lot of interesting work remains to be done …

# Future Work

♦ **Improving efficiency** of learning empirical explanations.

– Techniques presented are embarrassingly parallel.

♦ Use **supervised learning** for empirical explanations.

– Current techniques use unsupervised techniques.

♦ **Combined learning** of constraints and empirical explanations.

– Constraints used for data quality tend to be relatively simple.

– Empirical explanations can be more complex.