

Data Polygamy: The Many-Many Relationships among Urban Spatio-Temporal Data Sets

Fernando Chirigati¹, Harish Doraiswamy¹, Theodoros Damoulas^{2,3}, Juliana Freire¹



¹New York University

²The University of Warwick

³Alan Turing Institute



Urban Data Sets are Polygamous!

There are multiple interactions between entities of a city. These are captured by the **relationships** between urban data sets.

Relationship Queries

Find all data sets *related* to a given data set **D**

Enable *hypothesis generation* and *hypothesis testing*!

Hypothesis Testing

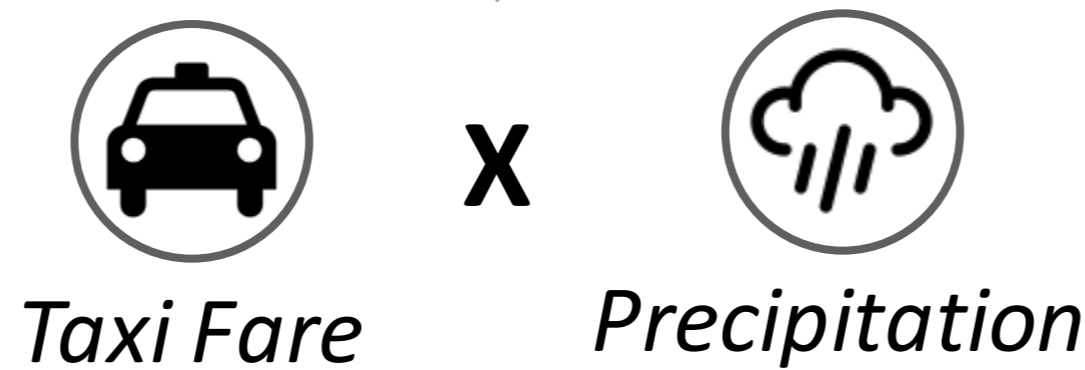
NYC residents often struggle to get a taxi when it is raining.

Long-standing hypothesis:

- Taxi drivers set an income goal
- They reach goal faster on rainy days

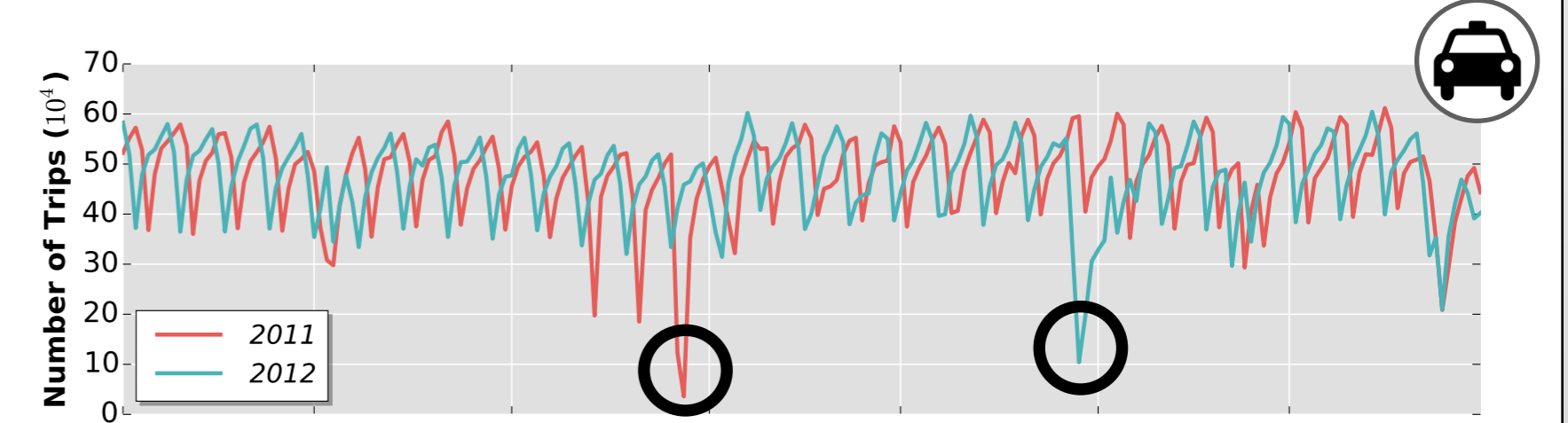
Can we test such hypothesis? **Yes!**

Find relationships between **Taxi and Weather** data sets



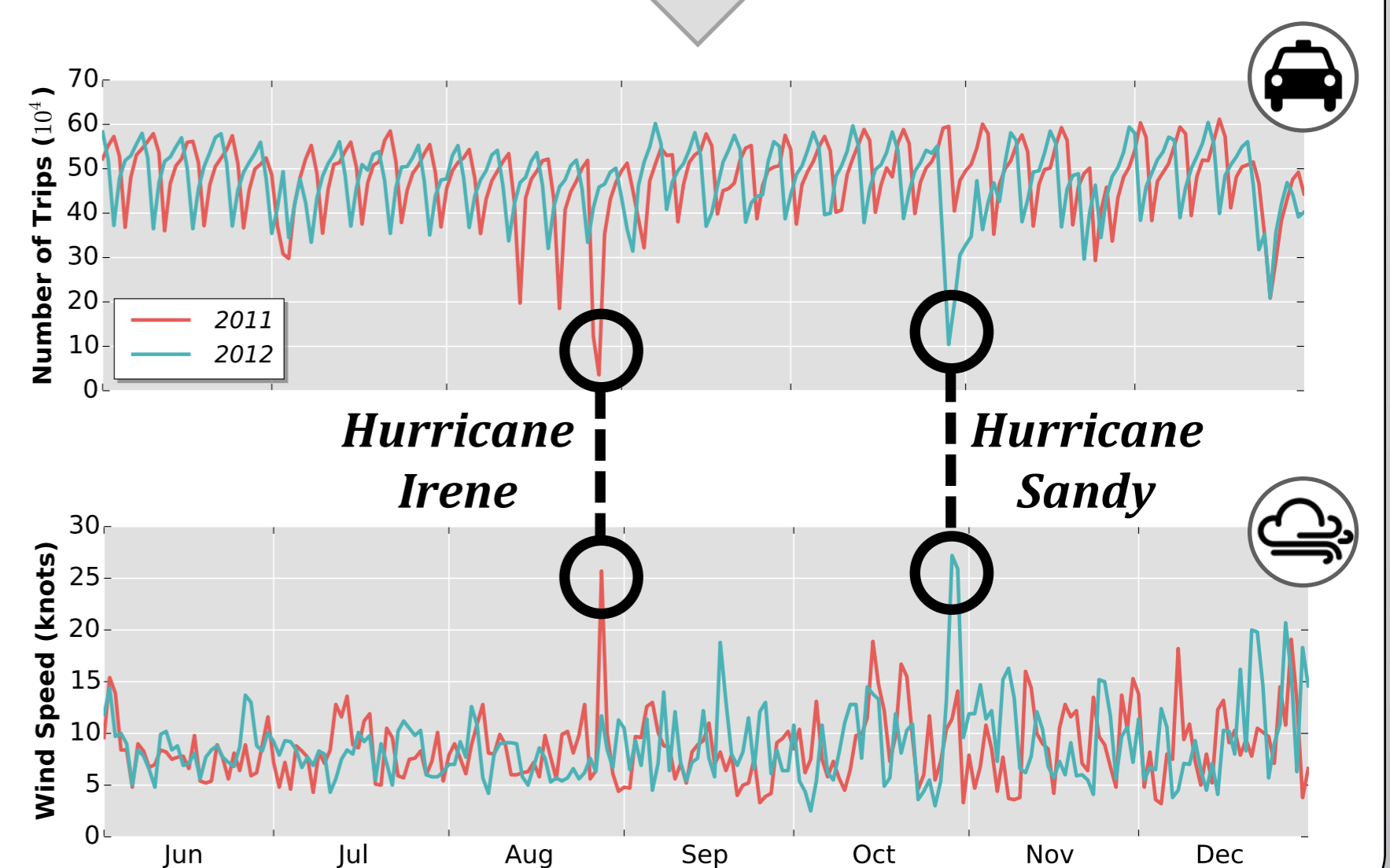
Hypothesis Generation

While analyzing the NYC Taxi data set...



How to **explain** these features?

Find all data sets related to the **Taxi** data set

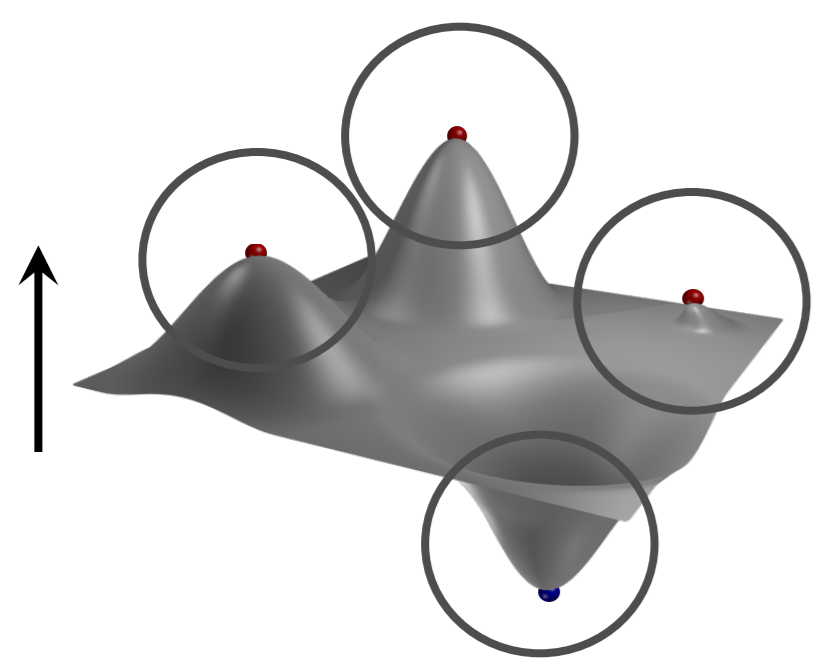


Challenge 1: How to define a data set relationship?

Our Approach: Computational Topology

1) Modeling the Data as a Terrain

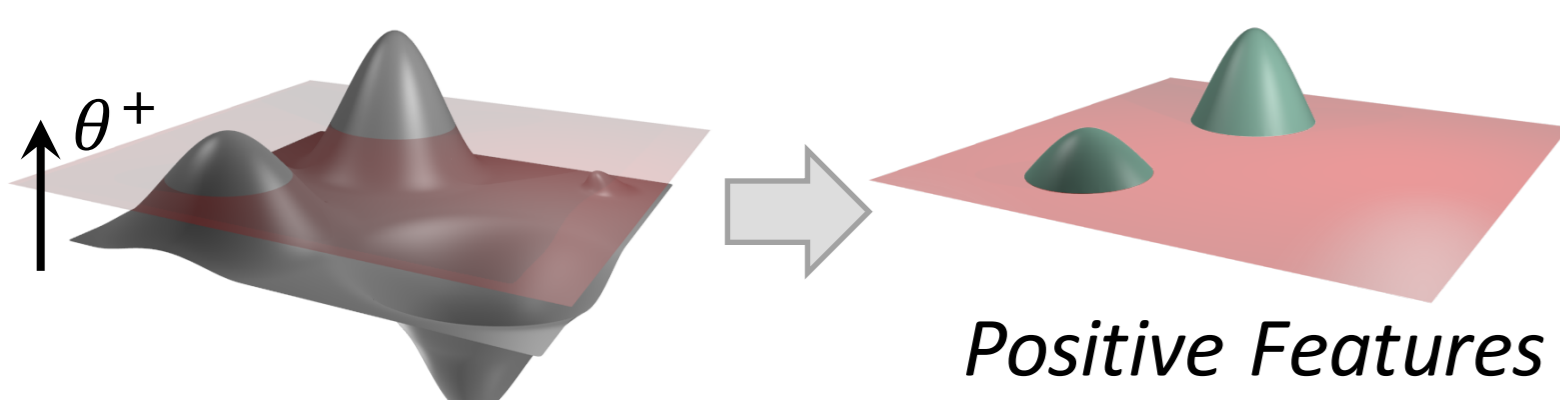
$$f : [S \times T] \rightarrow \mathbb{R}$$



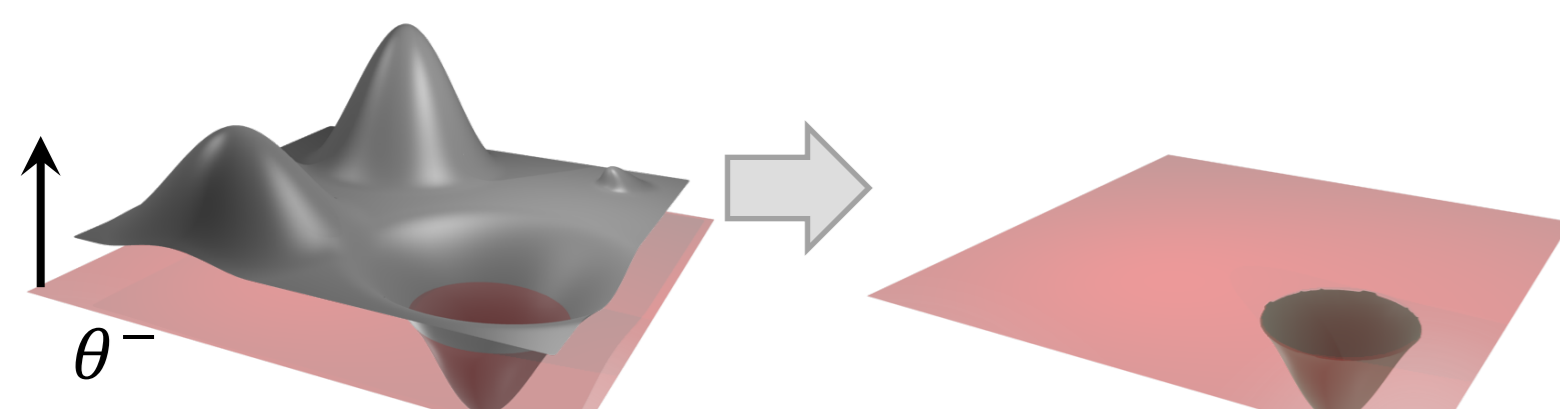
- Critical points
- Peaks
- Valleys

2) Identifying and Computing Topological Features

Neighborhoods of critical points

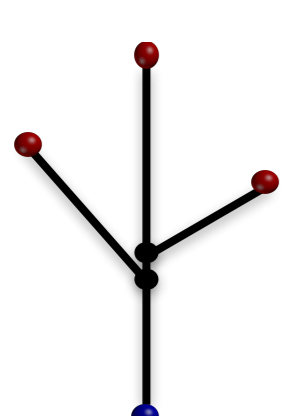


Positive Features



Negative Features

Index: Merge Tree

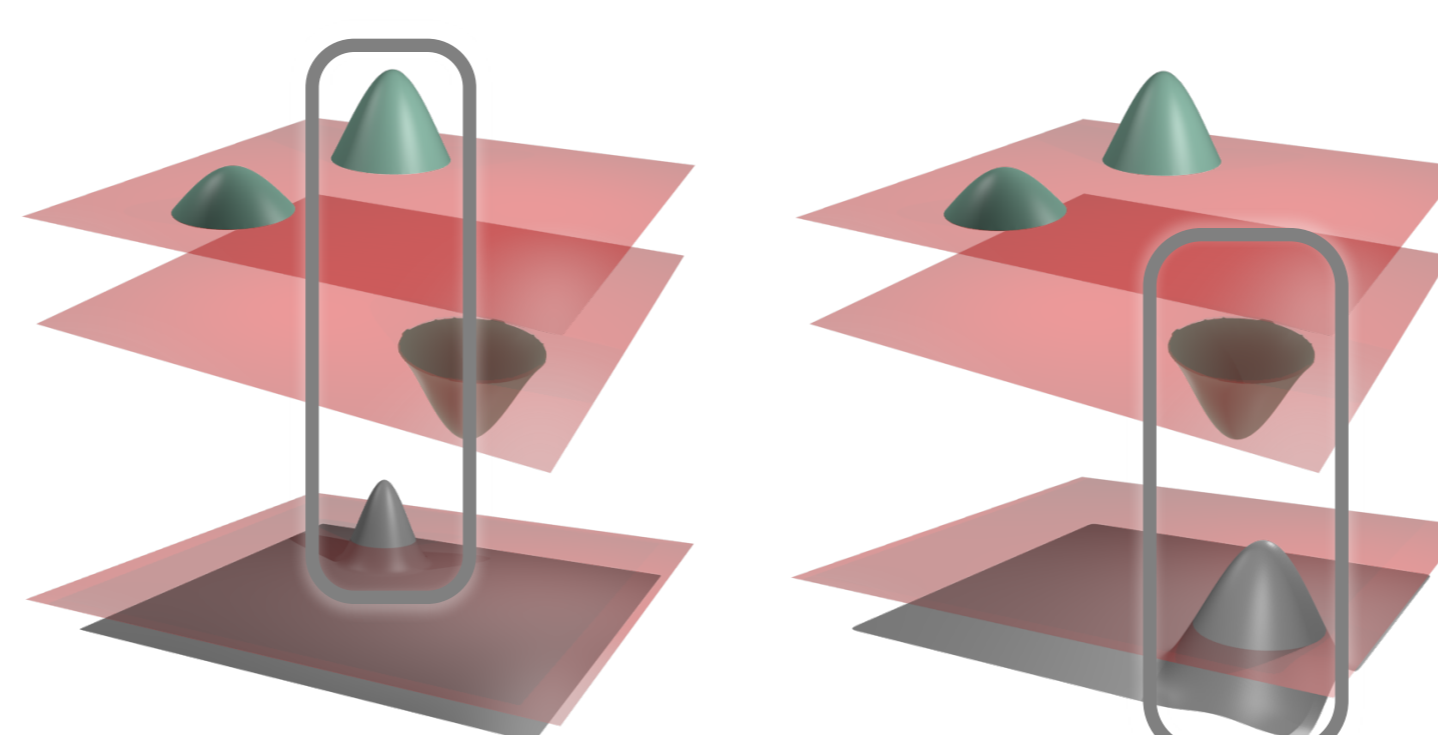


Thresholds θ^+ and θ^- computed in a data-driven approach

Computing features is **output sensitive!**

3) Identifying Topology-based Relationships

Relationship between features



Positive Relationship

Negative Relationship

Relationship between functions

- Relationship Score (τ)
Nature of the relationship

$$\tau = \frac{\#p - \#n}{|\Sigma|}$$

- Relationship Strength (ρ)
How often they are related

$$precision = \frac{\#tp}{\#tp + \#fp} \quad recall = \frac{\#tp}{\#tp + \#fn}$$

$$\rho = F_1(f_1, f_2) = 2 \times \frac{precision \times recall}{precision + recall}$$

Challenge 2: Data Complexity

- Multiple spatio-temporal resolutions
- Large data sets
- Relationships can be between any of the attributes

meaningful relationships

needle in a haystack

Our Approach:

- Monte Carlo tests filter potentially coincidental relationships
- Further filtering using τ and ρ

Reduces the number of output relationships in around 99%

Interesting Relationships

Taxi and Wind Speed

- No. taxis \times Wind speed

Taxi and Rainfall

- No. taxis \times Precipitation

+ Taxi fare \times Precipitation

Weather and Citi Bike

+ Snow precipitation \times Trip duration

- Snow precipitation \times Active stations

Weather is the most polygamous data set !