

#### **Scorpion**, a detective that is good at explaining the anomalies

#### Scorpion, explaining <u>outliers in aggregate queries</u>

• Aggregation function: SUM, AVERAGE, STD, MIN, MAX

# outliers in aggregate queries



Time	SensorID	Voltage	Humidity	Temp.
11AM	1	2.64	0.4	34
11AM	2	2.65	0.5	35
11AM	3	2.63	0.4	35
12PM	1	2.7	0.3	35
12PM	2	2.7	0.5	35
12PM	3	2.3	0.4	100
1PM	1	2.7	0.3	35
1PM	2	2.7	0.5	35
1PM	3	2.3	0.5	80

A table that records the attributes of sensors



# outliers in aggregate queries

Dataset from A hospital:

A table with one row per patient visit, and 45 columns that describe patient demographics, diagnoses, and other attributes describing the visit

SELECT **SUM**(expense), disease FROM hospital GROUP BY disease

lung cancer cases disproportionately account for millions of dollars

WHY?

### Scorpion :

### an interactive system answer "why" questions in the context of SQL aggregation queries.

Time	SensorID	Voltage	Humidity	Temp.	]			
HAM	1	2.64	0.4	34	j	Time	AVG(temp)	Label
11AM	2	2.65	0.5	35	] :	LIAM	34.6	Hold-out
HAM	3	2.63	0.4	35		1204	56.6	Outline
12PM	1	2.7	0.3	35		(D)(	50	Outling
12PM	2	2.7	0.5	35	1	IFM		Outlier
12PM	3	2.3	0.4	100		/		
IPM	1	2.7	0.3	35	/			
1PM	2	2.7	0.5	35				
IPM	3	2.3	0.5	80	1			

individual tuples are **less informative**, not enough for understanding why

Scorpion tells why:



predicate

- Predicates give broader, coarse-grained explanations
- Predicates: Conjunction of discrete subsets and continuous range



### Conjunction

Sensor Sensor &Voltage Sensor &Voltage & Light exponential in # of attributes

### Discrete

Sensor=1 Sensor= 1 or 2 Sensor=2 or 3 exponential in # of unique values

### Continuous

Sensor=1 Sensor= 1 or 2 Sensor=2 or 3 Quadratic in # of unique values

### Predicate space is very large

# scoring function

- We need a scoring function for ranking and returning to p explanations in the exponential space Scoring Function: Influence(P)
- Quantify how much predicate P influenced the result of aggregation function (sum, average)

Change in output

 $infl_{agg}(p) = \frac{1}{(\# o)}$ 

(# of records to make the change)

#### Change in output

## infl<sub>agg</sub>(p) =

### (# of records to make the change)

1	12PM	1	2.7	0.3	35
t	12PM	2	2.7	0.5	35
İ	12PM	3	2.3	0.4	100

Predicate1: sensor=3 One tuple causes the change Predicate2: sensor=3 or 2 Two tuples cause the change infl(P2)=21.6/2=10.8 Predicate3: sensor=3 or 2 or 1 Three tuples cause the change infl(P3)=51.6/3=17.2

infl(P1)=(56.5-35)/1=21.6

#### Change in output infl<sub>agg</sub>(p) (# of records to make the change)

#### Lamda is a hyperparameter : Leave the choice to the user

#### Lamda=1

Predicate1: sensor=3

Predicate2: sensor=3 or 2 One tuple causes the change Two tuples cause the change

#### lamda=0

Predicate3: sensor=3 or 2 or 1 Three tuples cause the change

The higher the lamda, the more selective predicate it produces

# Change in output 🗙 🔰

### infl<sub>agg</sub>(p) =

(# of records to make the change)

V: users indicate whether the outliers are too high or too low

Time	AVG(temp)	Label	v
11AM	34.6	Hold-out	
12PM	56.6	Outlier	<+1>
1PM	50	Outlier	< +1 >

	~			~~
12PM	1	2.7	0.3	35
12PM	2	2.7	0.5	35
12PM	3	2.3	0.4	100

Predicate1: sensor=3 Infl(P1)=26.1\*1=26.1 Predicate2:sensor=1 Infl(P2)=-10.9\*1=-10.9 Predicate1: sensor=3 Infl(P1)=26.1\*(-1)=-26.1 Predicate2:sensor=1 Infl(P2)=-10.9\*(-1)=10.9

$$inf_{agg}(o, h, p, v_o) = \bigcap inf_{agg}(o, p, v_o) - (1 - \bigcap |inf_{agg}(h, p)|$$

C: users indicate how much the predicate should take the hold- out results into consideration

Time	AVG(temp)	Label	V
11AM	34.6	Hold-out	-
12PM	56.6	Outlier	< +1 >
1PM	50	Outlier	< 11 >

$$inf_{agg}(O, H, p, V) = \bigcap_{o \in O} \frac{1}{|O|} \sum_{o \in O} inf_{agg}(o, p, v_o) - (1 - \bigcap_{h \in H} \max |inf_{agg}(h, p)|$$

The user often select multiple outliers results and hold-out results, we extend the notion by averaging the influence over the outlier results

Time	AVG(temp)	Label	v
11AM	34.6	Hold-out	-
12PM	56.6	Outlier	<+1>
1PM	50	Outlier	< +1 >

### **Influential Predicates Problem**

$$p^* = \arg \max_{p \in P_{A_{rest}}} \inf(p)$$

# **Scorpion Architecture**



# A naïve implementation





Explore some properties of aggrega — — > Enable more efficient implementations

### some properties

Incrementally removable





Influ(P)=15-SUM({1,2,3})

Influ(P)=15-(15-SUM({4,5}))=SUM({4,5})

SUM, AVG, STD, COUNT, MAX, MIN

### some properties

- Incrementally removable
- Independent

Scorpion  

$$p^* = \arg \max_{p \in P_{A_{rest}}} inf(p)$$

- 1.the influence of a set of tuples strictly depends on the influen ces of the indi vidual tuples
- 2. adding a tuple t more influential than t \*  $\in$  T with minimum influence ca n't decrease the set's influence



# Least Most influence



**TOP-DOWN** decision tree-based algorithm that recursively partitions the predicates and merges similar(adjacent) predicates



## some properties

- Incrementally removable
- Independent
- Anti-monotonic



Non-influential tuple sets can not contain influential sub-tuple sets: prune the search space!



#### **bottom-up**:

first search for influential single-attribute predicates, then intersect them to construct multi-attribute predicates.







# Conclusions

- 1st system to explain outliers in aggregation quries
- general system for a large class of why questions
- optimization based on operator properties rather than for specific scenario
- clean data automatically and interactively

Thank you! Q&A