

# Outlier detection

From statistical analysis to  
more advanced ML/DL  
approaches

**Raphael Espanha**  
Data Scientist

October 15th, 2019

**Know** the unknown.

**DSPT**  
DATA SCIENCE PORTUGAL

**W e D O**  
technologies

A MOBILEUM Company

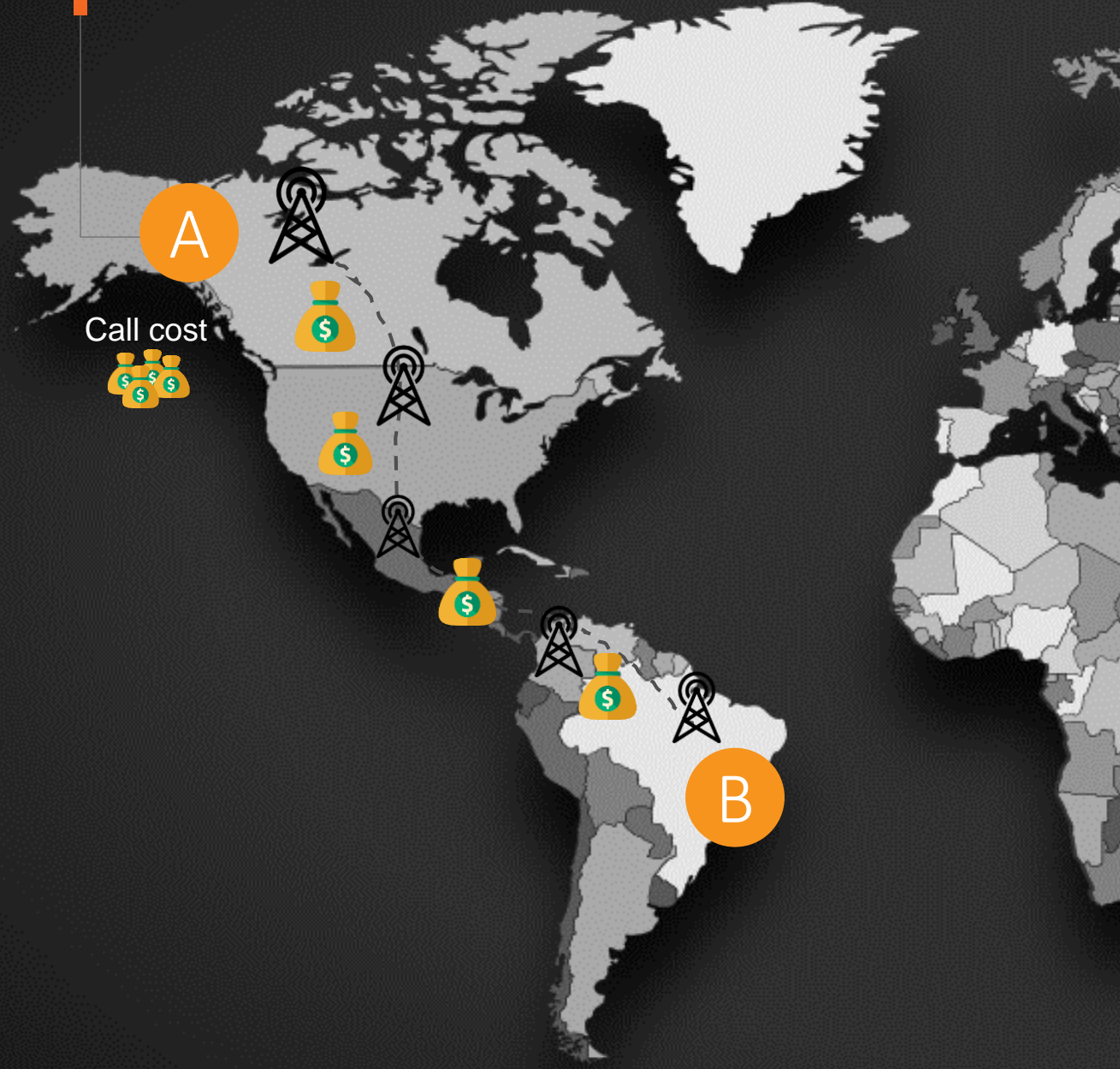


# BYPASS

MTR USE CASE

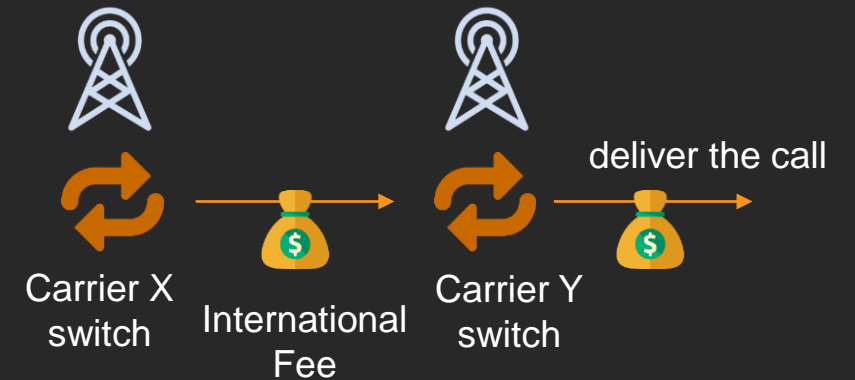


# International Call



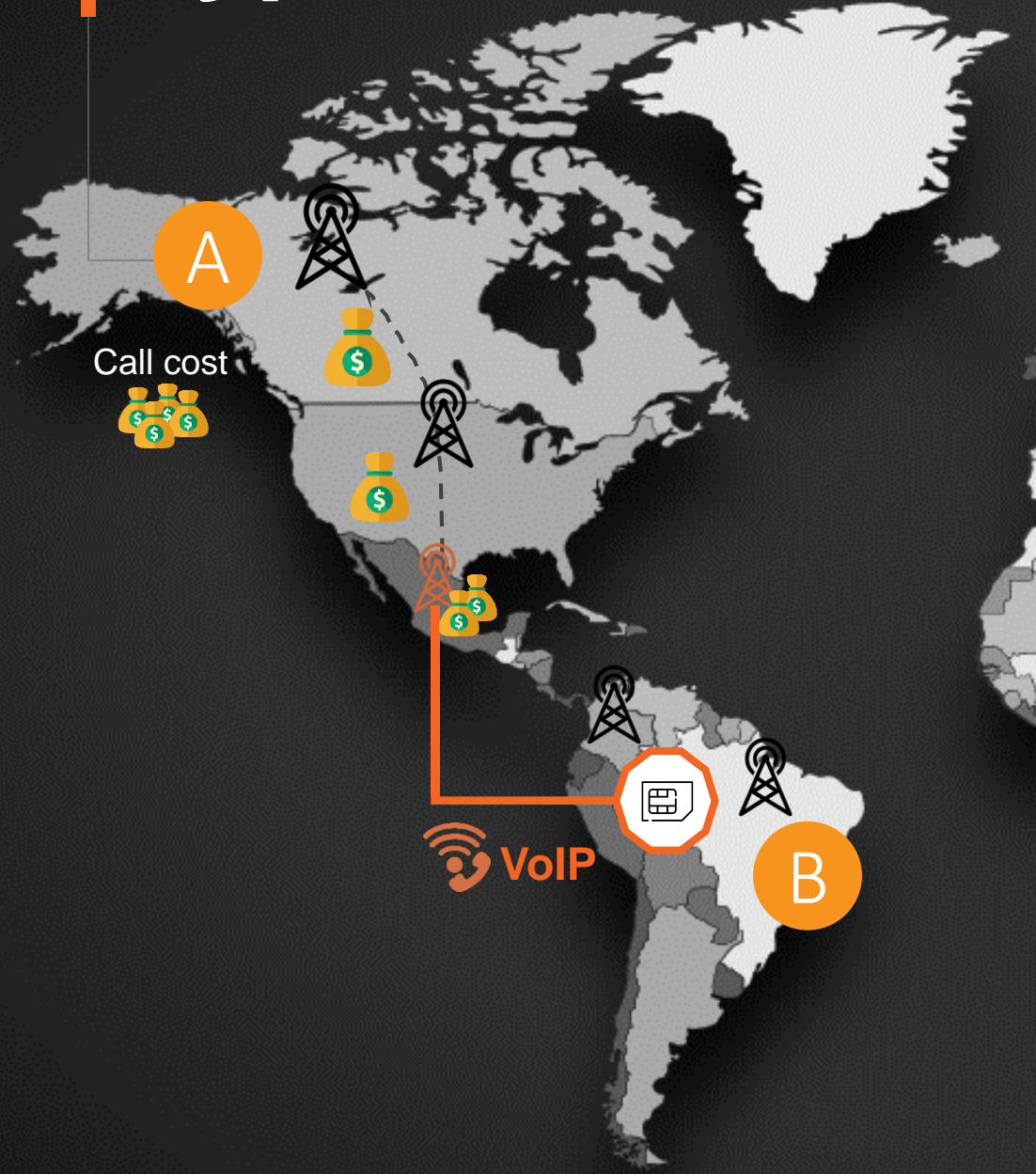
Phone companies have agreements with each other and use the services of carriers

By paying interconnection fees



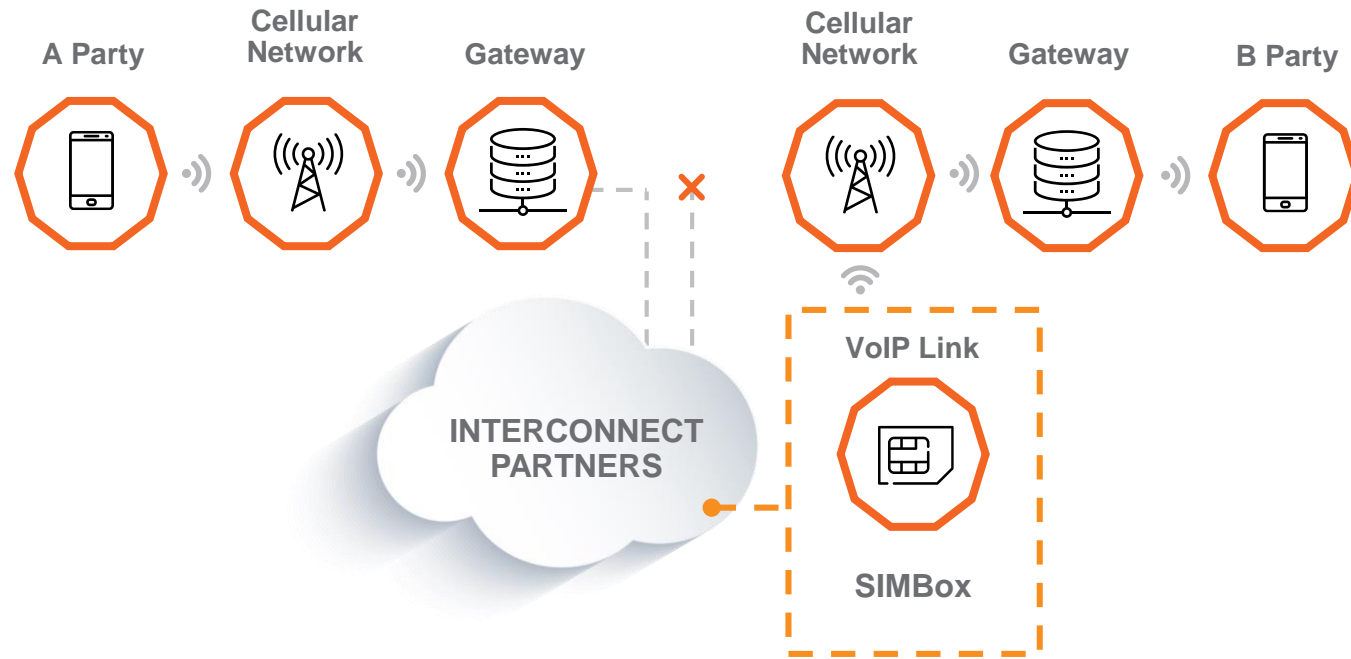


# Bypass



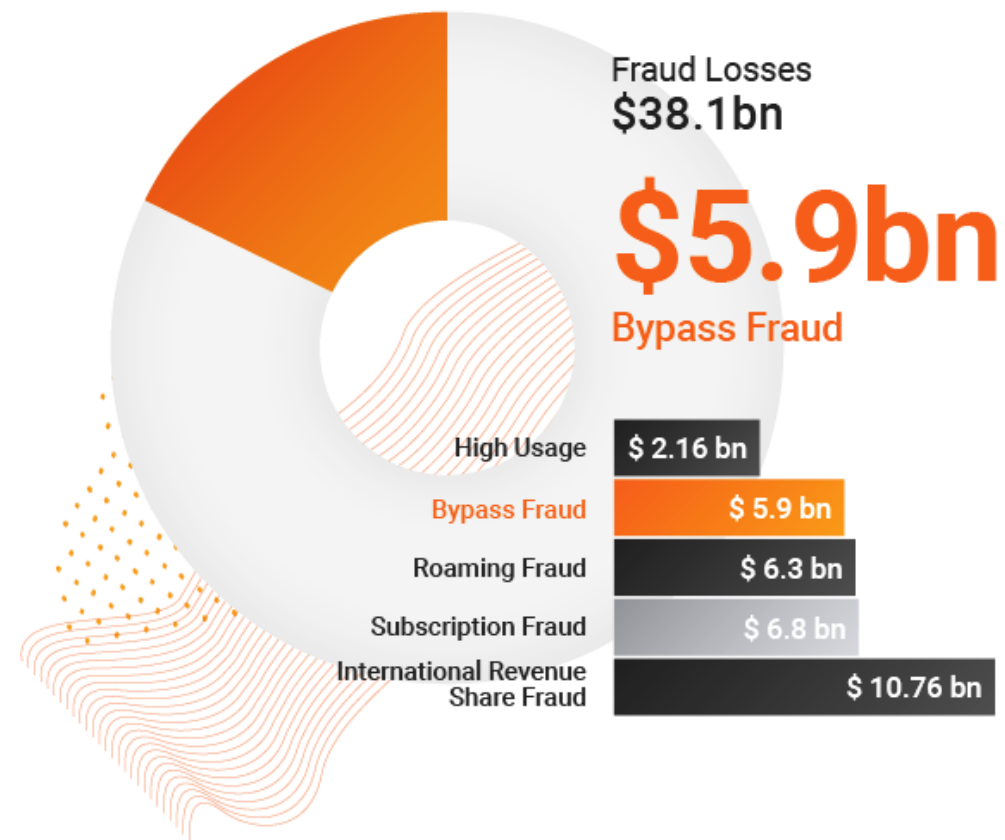
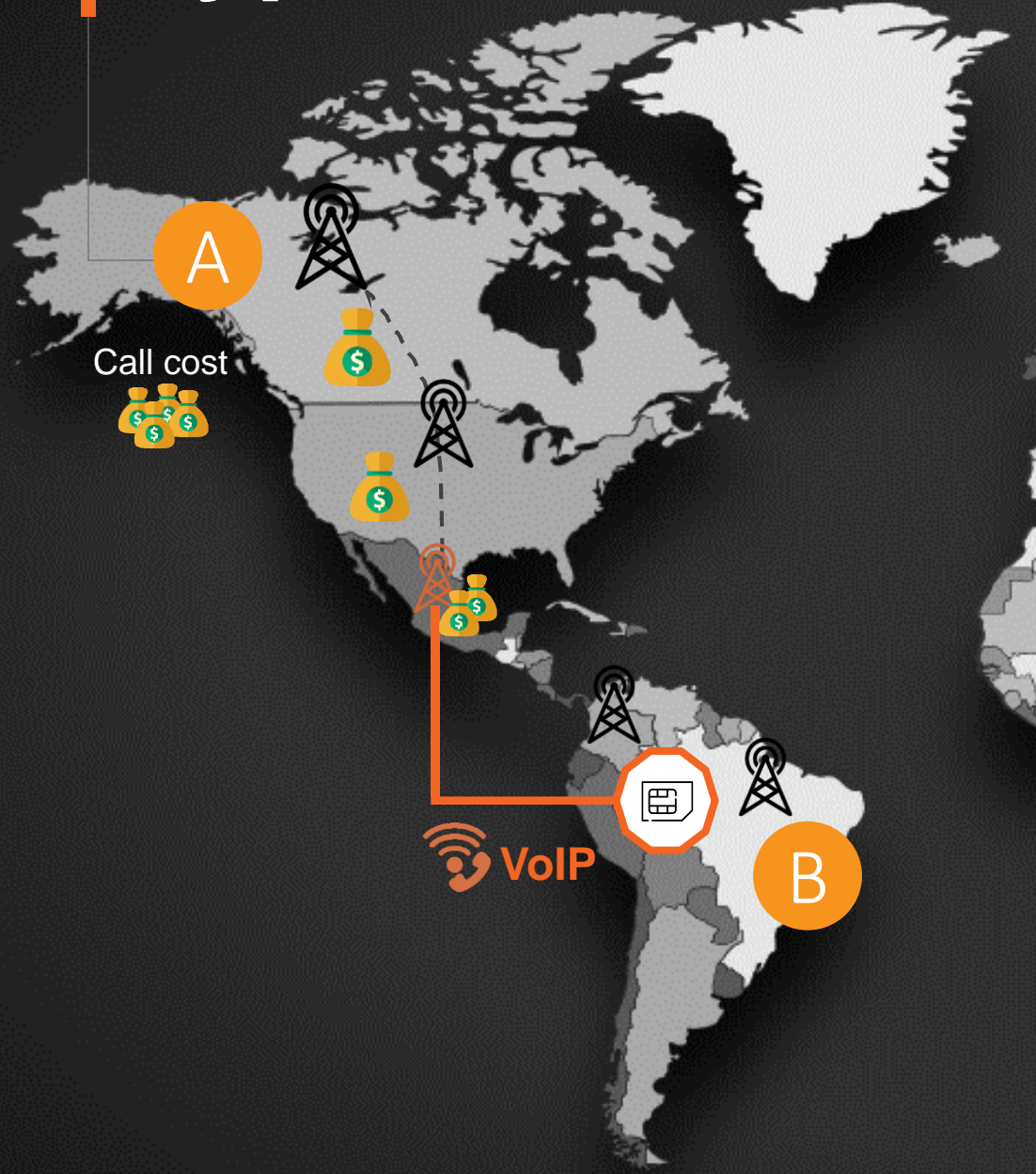
COUNTRY A (Mexico)

COUNTRY B (Brazil)



$$\text{Fraudster profit} = \text{International Fee to forward the call} - \text{Local call cost}$$

# Bypass



## Financial Impact

- Loss of revenue for legitimate operators

## Image Impact

- No call completion
- Odd or no calling number
- Bad call quality



**GOAL**



# Goal – Detect abnormal ranges

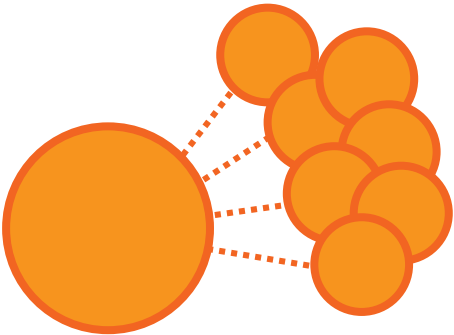
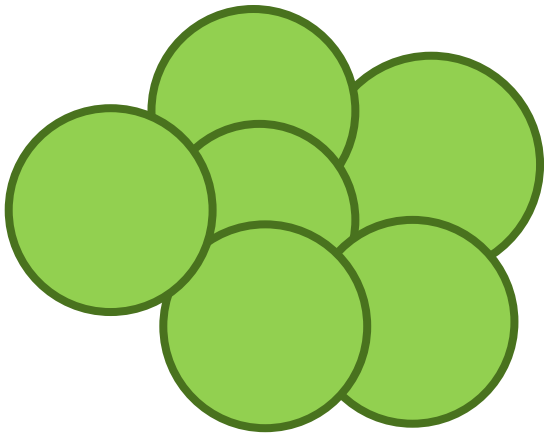


a_number
ABCDE
ABCDF
ABCDG



range_1
ABCD

Normal ranges



A numbers

Range with abnormal behavior

**DATA**

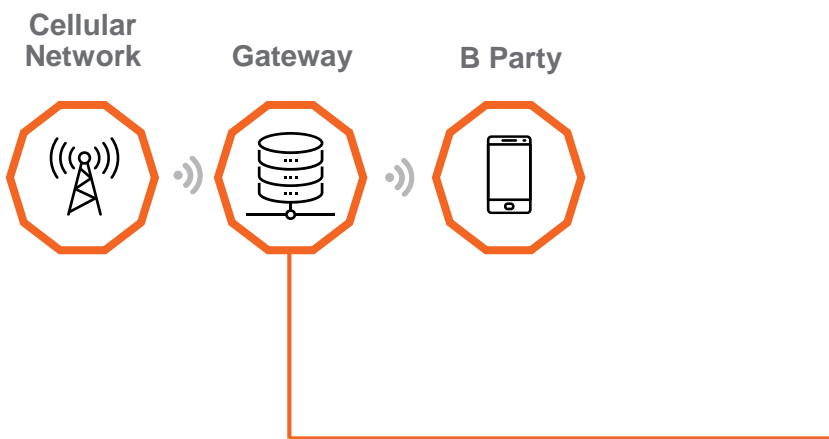




# Data – Signaling events

Signaling data from august 2018

COUNTRY B (Brazil)



a\_number - origin number  
b\_number - destiny number  
time: call timestamp

a_number ↕	b_number ↕	time ↕
ZNFNNUULP	BNFVZNOPOIUFP	20180806211322
ZIIPZIUZ	BNFVZNFROIRFO	20180816191720
IIFLIZPUOFL	BNFVZNLFFVNOIO	20180820154021
LLRRNRUIUZUI	BNFVZNNIUPZVP	20180814122533
ZNFZUPPPR	BURVDUZNZUUFOLIFNFP	20180815204314
VNNVVFRLZ	BNFVZNOLOUVNV	20180806160006
IOVFFNVOOLNF	BURVZNNIUFRUV	20180804105707
LLRIZIILIPFR	BNFVZNNIUZIUIO	20180814012423
ZNUPOZROO	BNFVZNNIUFOPL	20180813171510
ZNOZNULOZ	BNFVZNVOLLLNO	20180815141804

# Data Cleansing

Client requested filters were applied:

- A numbers from starting with L or I
- A number length > 9
- Null values

a_number	b_number	time
ZNFNNUULP	BNFVZNOPOIUFP	20180806211322
ZIIPZIUZ	BNFVZNFROIRFO	20180816191720
IIFLIZPUOFL	BNFVZNLVFNIO	20180820154021
LLRRNRUIUZUI	BNFVZNNIUPZVP	20180814122533
ZNFZUPPPR	BURVDUZNZUUFOLIFNFP	20180815204314
VNNVVFRLZ	BNFVZNOLOUVNV	20180806160006
IOVFFNVOOLNF	BURVZNNIUFUVR	20180804105707
LLRIZIILIPFR	BNFVZNNIUIZUO	20180814012423
ZNUPOZROO	BNFVZNNIUFOPL	20180813171510
ZNOZNULOZ	BNFVZNVOLLLNO	20180815141804

33 929 956 events



a_number	b_number	time
IIFVNZFFNZR	BNFVZNNIUZIZL	20180821201213
IOVFZNPZZLUZ	BURVZNNIUFUVR	20180805114643
ILFNFVNPURP	BNFVZNNIUOILR	20180824190610
LLRPILRNPLUU	BNFVZNNIUZRNZ	20180820130935
IOVFNVRPZLU	BURVZNNIUFIPU	20180805202844
LVUFUOLNRFLZ	BNFVZNNIULIVI	20180802125215
LURFVPIVLVN	BURVZVULLUPRU	20180815202426
IPUZZUOIIRVP	BURVZVUIIVPNU	20180803114237
IIFNUNRLPZV	BNFVZNNIURNNO	20180806085534
IIFUIPPNVNU	BNFVZNNIURLIU	20180816120950

20 636 037 events

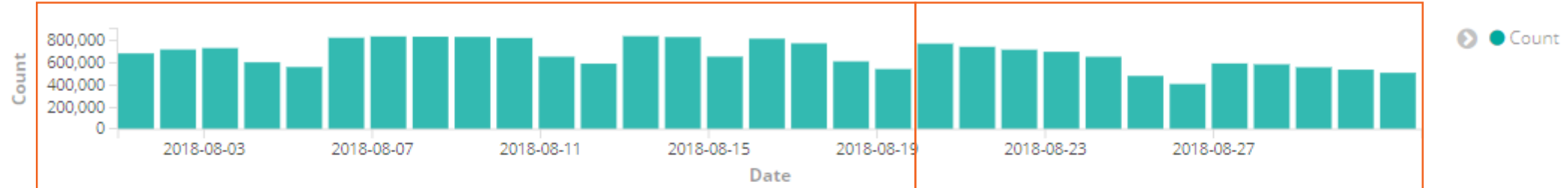


# Data Analysis – Time and distinct values

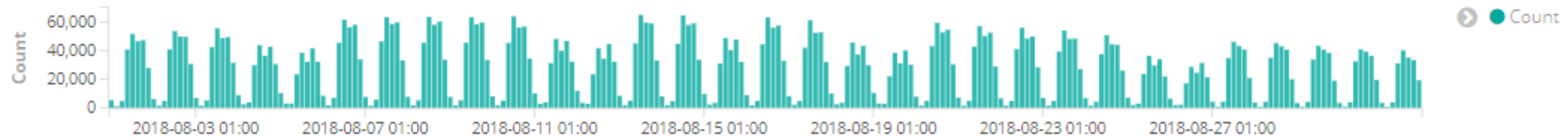
Signalling Daily Counts

Train

Test



Signalling Hourly Counts



Signalling Unique A Numbers

**4,445,792**

Unique A Numbers

Signalling Unique B Numbers

**1,325,114**

Unique B Numbers

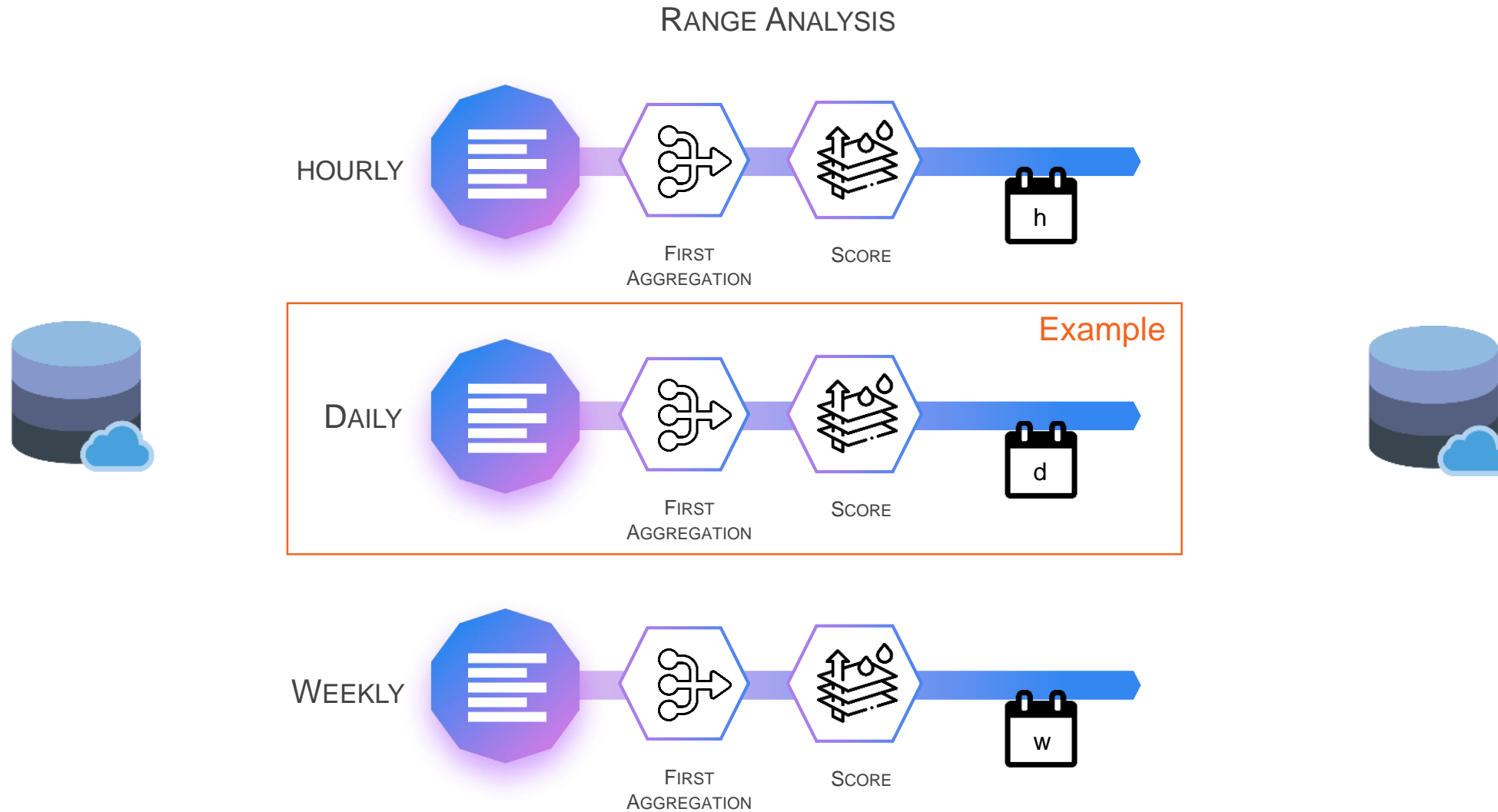
# Outlier Detection Range Numbers

Statistical approach

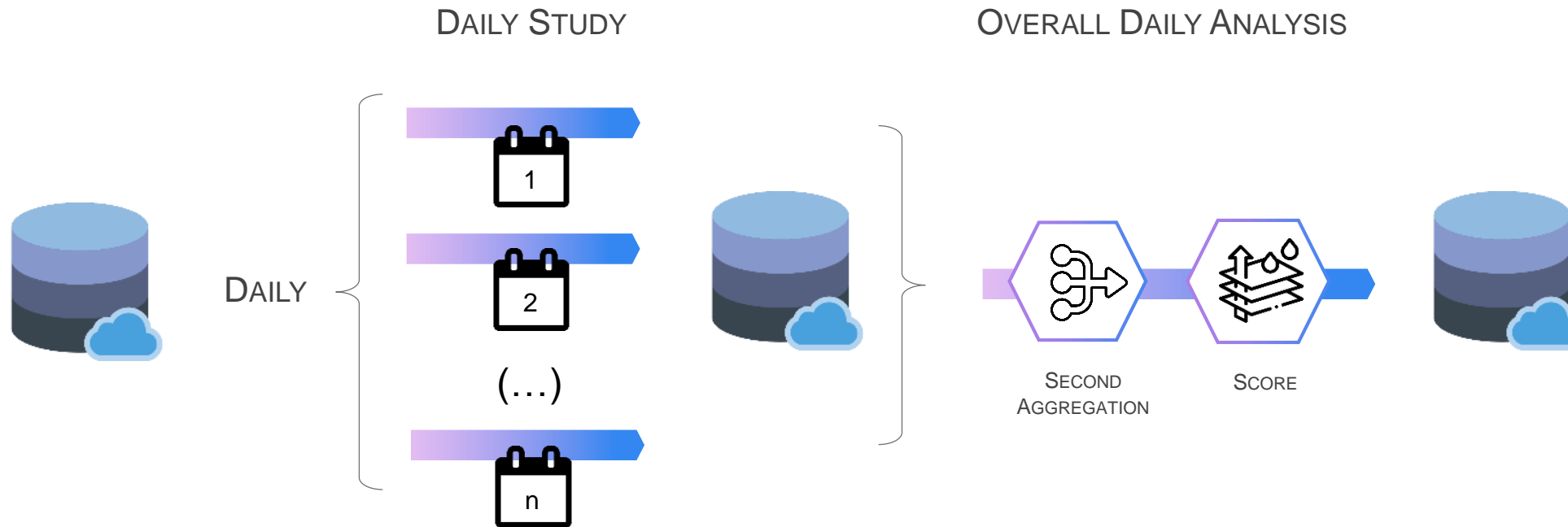




# Statistical Approach – Multiple time context analysis



# Statistical Approach – Daily analysis overview





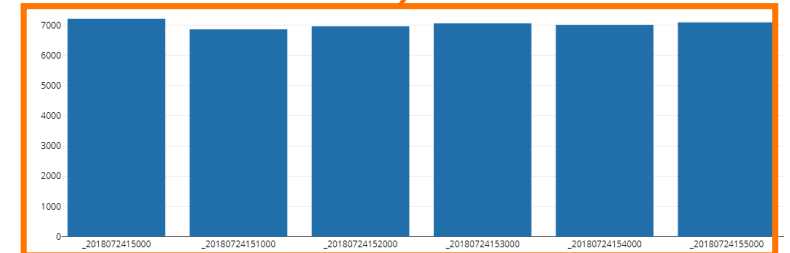
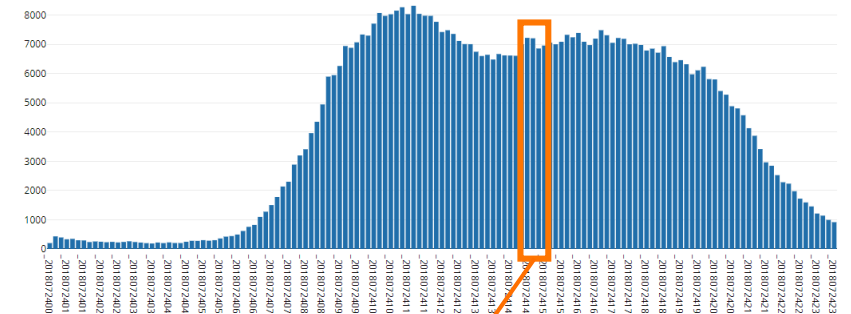
# Statistical Approach – Preprocessing - Time binning

Creating intervals of 10 min

10:43	10:40
10:45	10:40
10:47	10:40
10:51	10:50
10:55	10:50
10:59	10:50



time	interval_bin
20180819160732	201808191600
20180806114005	201808061140
20180824092934	201808240920
20180822155336	201808221550
20180807131051	201808071310
20180802085920	201808020850
20180820100427	201808201000
20180825164352	201808251640
20180812123831	201808121230
20180815173523	201808151730
20180824152437	201808241520

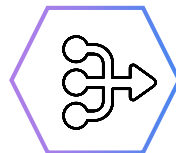


Equal-width Binning

# Statistical Approach - First aggregation (daily analysis example)

Group by a\_number

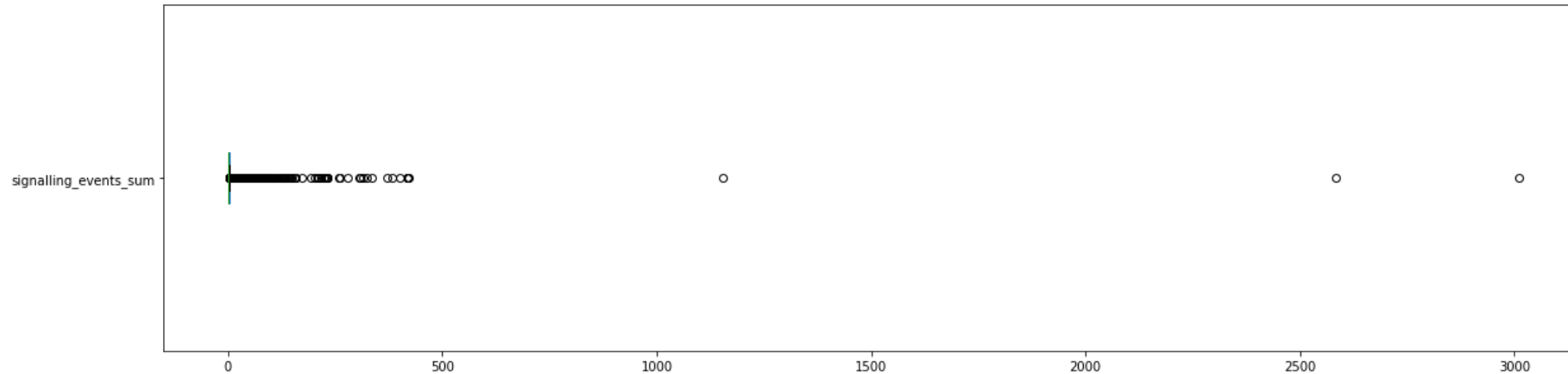
a_number	b_number	time	interval_bin
LLRPPNFVOZPL	BNFVZNNUFNFP	20180801143033	201808011430
IIFOUVZLLUO	BNFVZNNUILZUI	20180801182655	201808011820
IIFRRONLVNU	BNFVZNNUORRR	20180801193622	201808011930
LLROILIOPRIU	BURVZNLRNIZO	20180801133657	201808011330
INFIINRFOOU	BNFVZNNUZUPZ	20180801174707	201808011740
LZLNRNRZFU	BNFVZNNUOZVO	20180801115318	201808011150
LLROFPNLIZOO	BNFVZNNUPLZU	20180801190527	201808011900
IIFVRIFOUNF	BNFVZNNUFFVZ	20180801210050	201808012100
ILFUPUZORFR	BURVZNIPIONZV	20180801101755	201808011010
IIRFZZUFZNV	BNFVZNZIPPIVR	20180801123153	201808011230
IVLZNVIPVFF	BNFVZNNUPVVOV	20180801124409	201808011240
IIFNFVOIZZO	BURVZNZFVORU	20180801201836	201808012010
IIFUZZPPUFP	BNFVZNNUOLRU	20180801171554	201808011710



a_number	distinct_b_numbers	distinct_interval_bin	signalling_events
IIFZOVRUOI	1	1	1
IIFNOOZPZRN	1	1	1
LLVURNFFZUPI	1	1	1
LFRIURNFPIV	1	1	1
LLRPFRINZUVU	1	1	1
IIFVZVIZLZF	1	1	1
IIFFNNUVIIP	4	4	4
IZULVVPVIZU	1	1	1
LLVURPIINZII	1	1	1
ILFRZNRPRZR	1	1	1
LLRLZILUPZFP	2	1	2
IIFNZIVIVVN	1	1	1
LRZORIULII	1	1	1

# Statistical Approach – dropping low event ranges

Assuming that a numbers and ranges with very low daily events are not fraudsters



Using Tukey's fences

$$\underbrace{[Q_1 - k(Q_3 - Q_1), Q_3 + k(Q_3 - Q_1)]}_{\emptyset} \quad 5 \text{ (k=3)}$$

Removing a\_numbers with less than 5 events in one day:

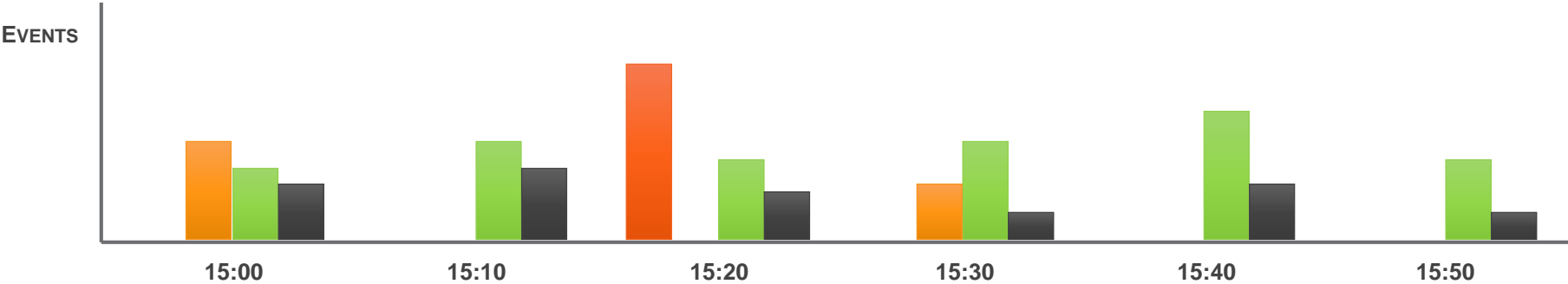
✓ Greatly improves algorithm performance

328 368 ranges to 15 490 (~4%)



# Statistical Approach – Event’s intensity

	a_number ↕	distinct_b_numbers ↕	distinct_interval_bins ↕	signalling_events ↕	event_intensity_score ↕
■	IFIULVNPLOR	16	1	20	0.9328
■	IFRUIRFUVIL	22	2	30	0.9036
■	IFRUZURNUIV	86	6	100	0.3165
■	LZVNNPPVFVFNZ	32	6	39	0.2633

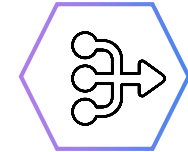


# Statistical Approach – Coefficient of Variation (CV)

Also known as relative standard deviation (RSD)

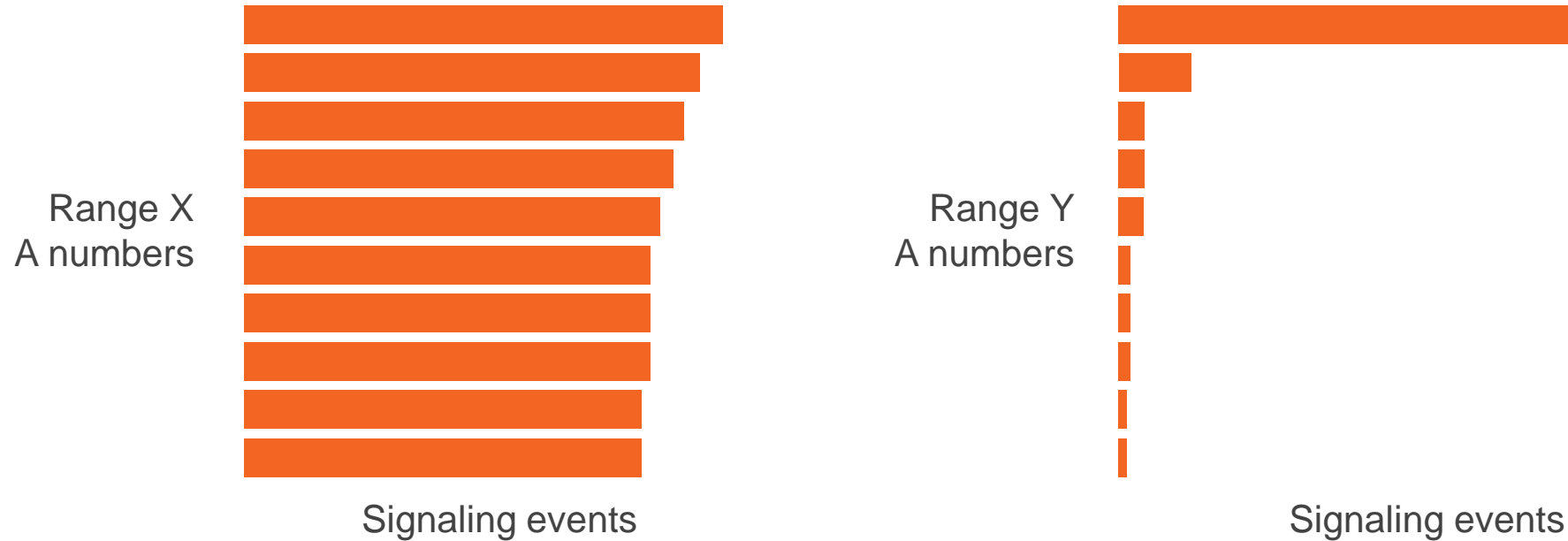
**Standardized measure** of dispersion of a **frequency distribution**

$$c_v = \frac{\sigma}{\mu}$$



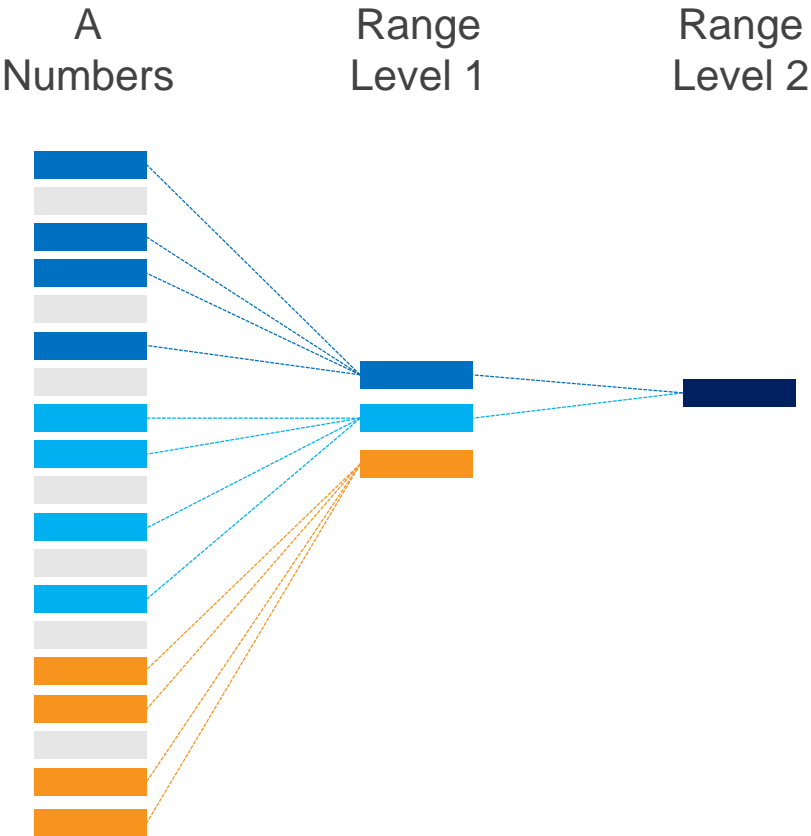
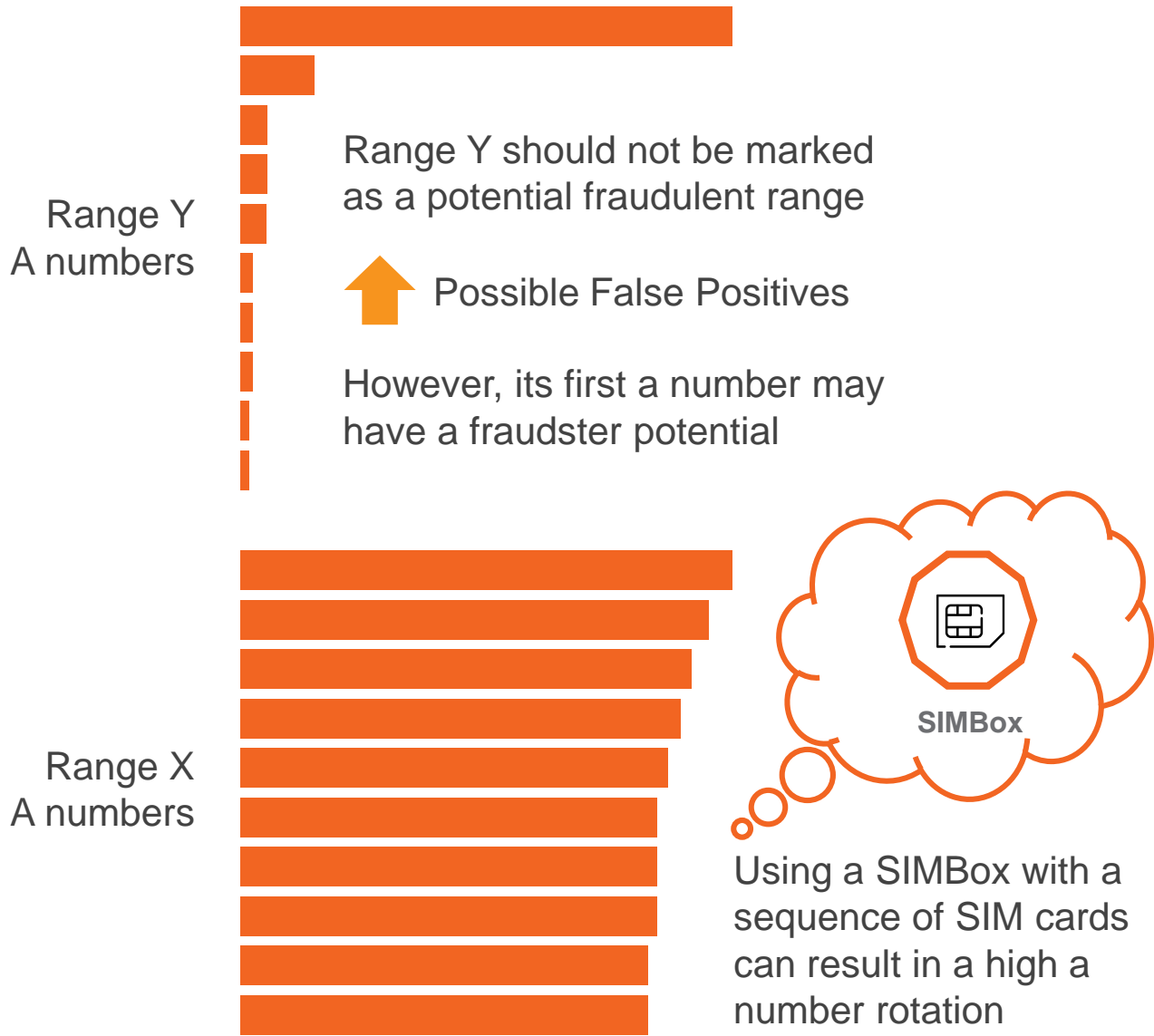
Requires group  
by range first

Compute a **rotation score** based on CV computed from events



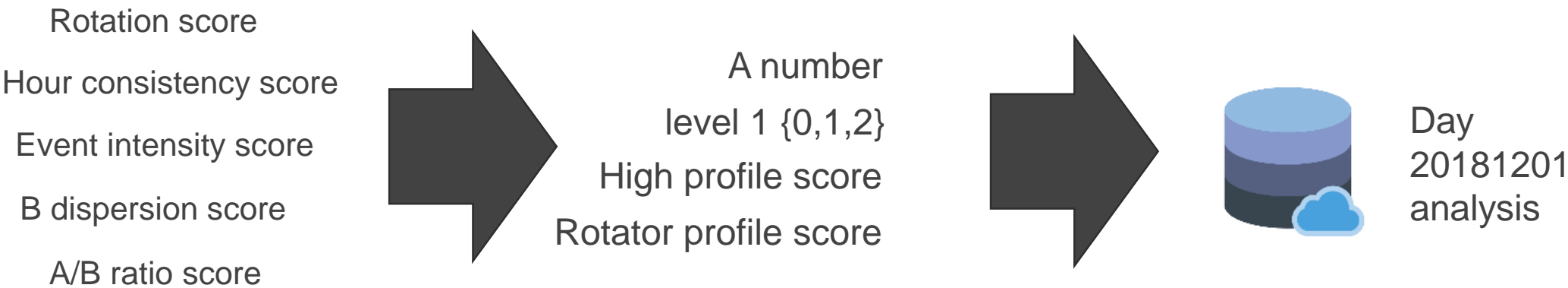
Range X rotation score >> Range Y rotation score

# Statistical Approach – Map back possible a numbers level up





# Statistical Approach – Overall Analysis and Profiles



Daily Analysis  
Finished

Overall Analysis

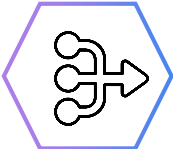


Day 20181201 analysis

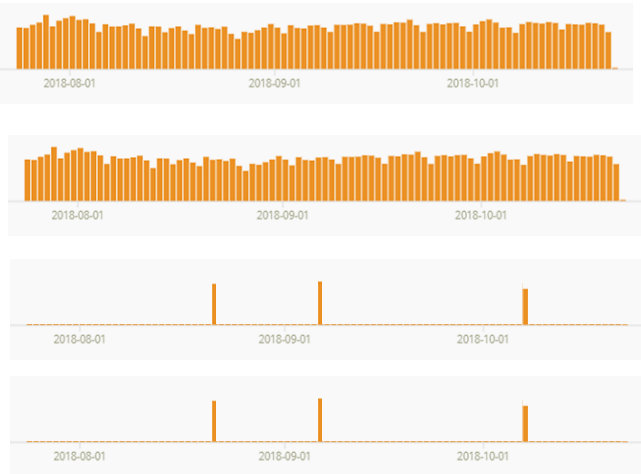
(...)



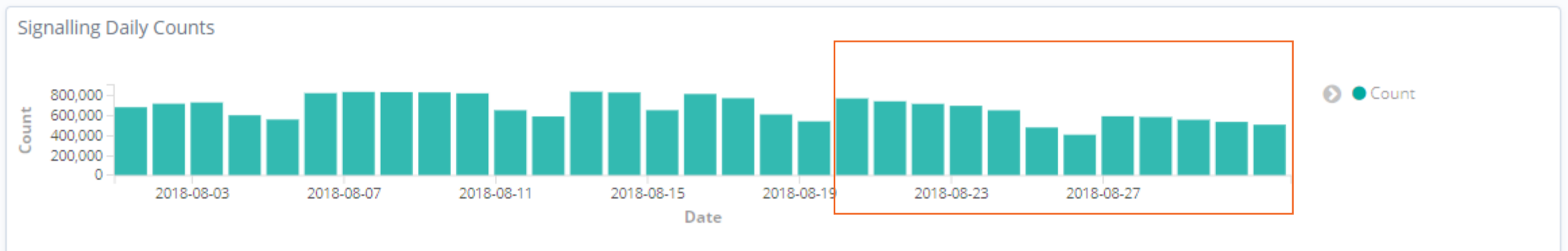
Day X analysis



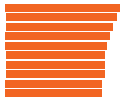
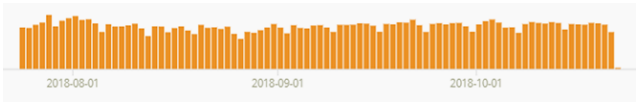
- Global Rotators
- Global High
- Local Rotators
- Local High



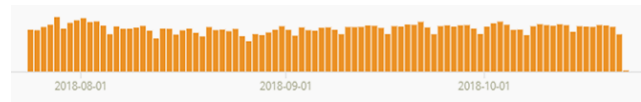
# Statistical Approach – Results



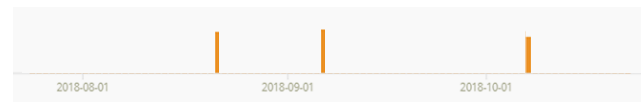
Global Rotators



Global High



Local Rotators



Local High



Tukey's fences  
K=3

383 cases

19 / 21 known bypass (90%)  
0 / 52 known FP (0%)  
365 unclassified

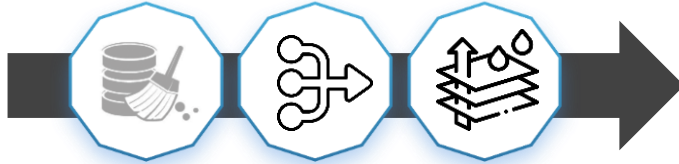
# Outlier Detection Range Numbers

Machine Learning - PCA





# Principal Component Analysis



Similar data preparation steps

- Range
- Range Level
- Signaling Events CV
- AB Ratio
- Rotation Score
- Interval consistency
- (...)

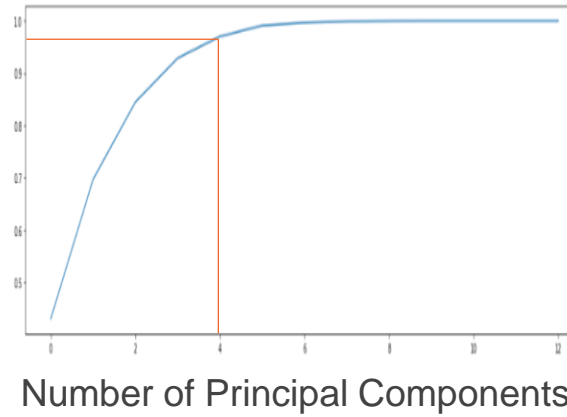
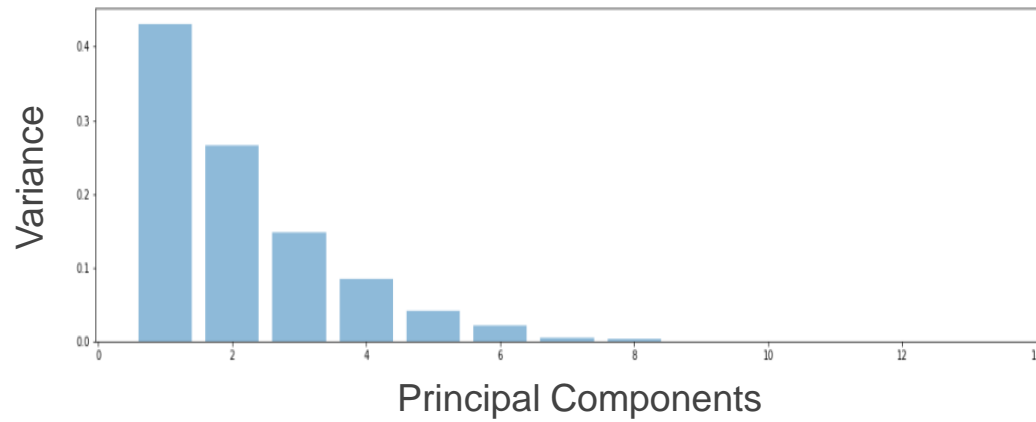


standardized scale (z)

$$z = \frac{x - \mu}{\sigma}$$

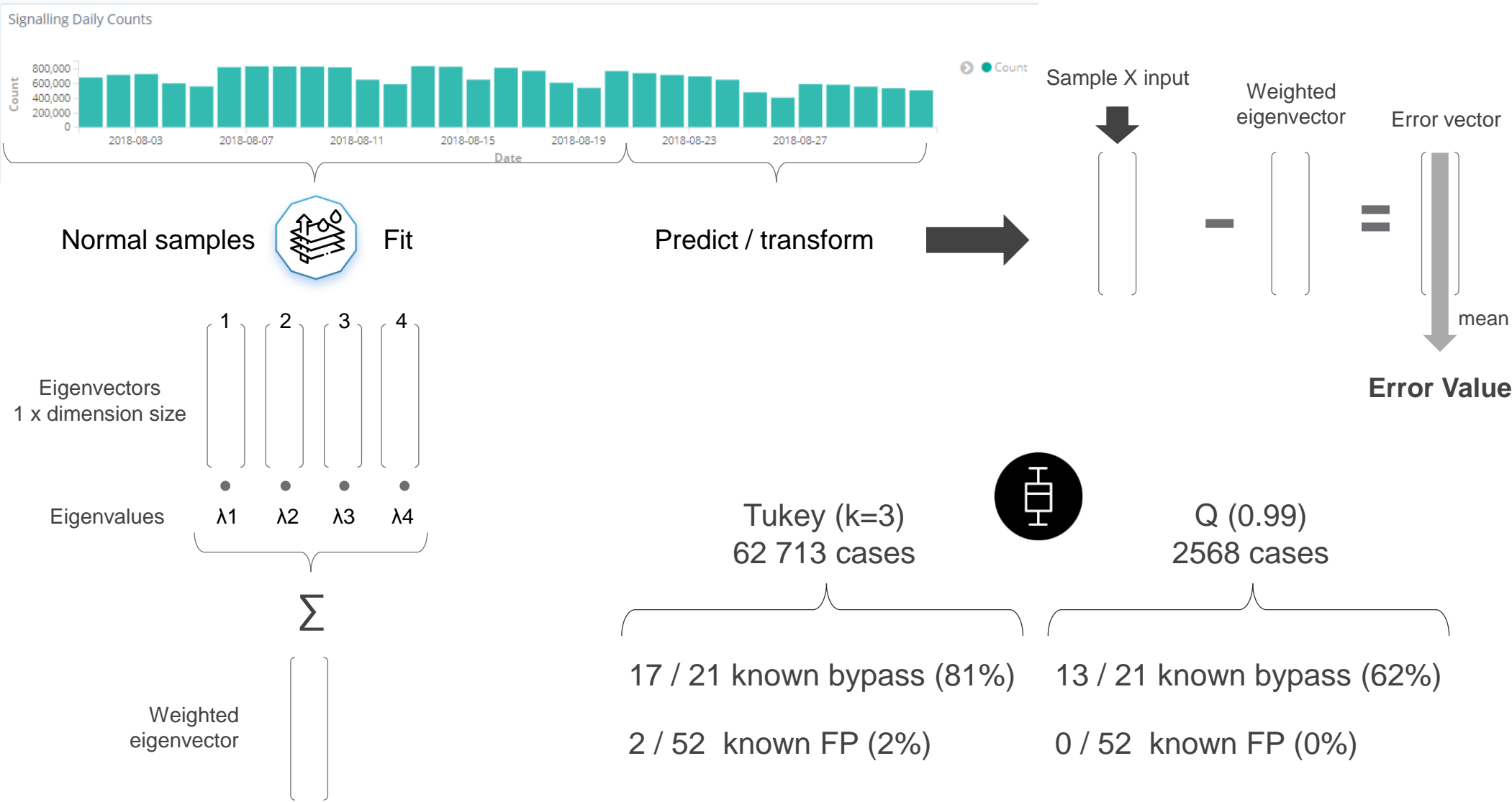
Principal Component Analysis is a **linear transformation** of data that reduces its dimension into **Principal Components (PCs)**

**Principal Components** are new variables that are constructed as linear combinations of the initial variable



PCs	Retained Information (%)
1	43
2	70
3	84
4	93

# Principal Component Analysis - Results

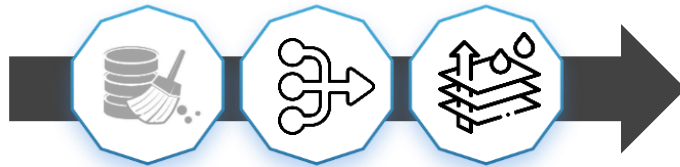
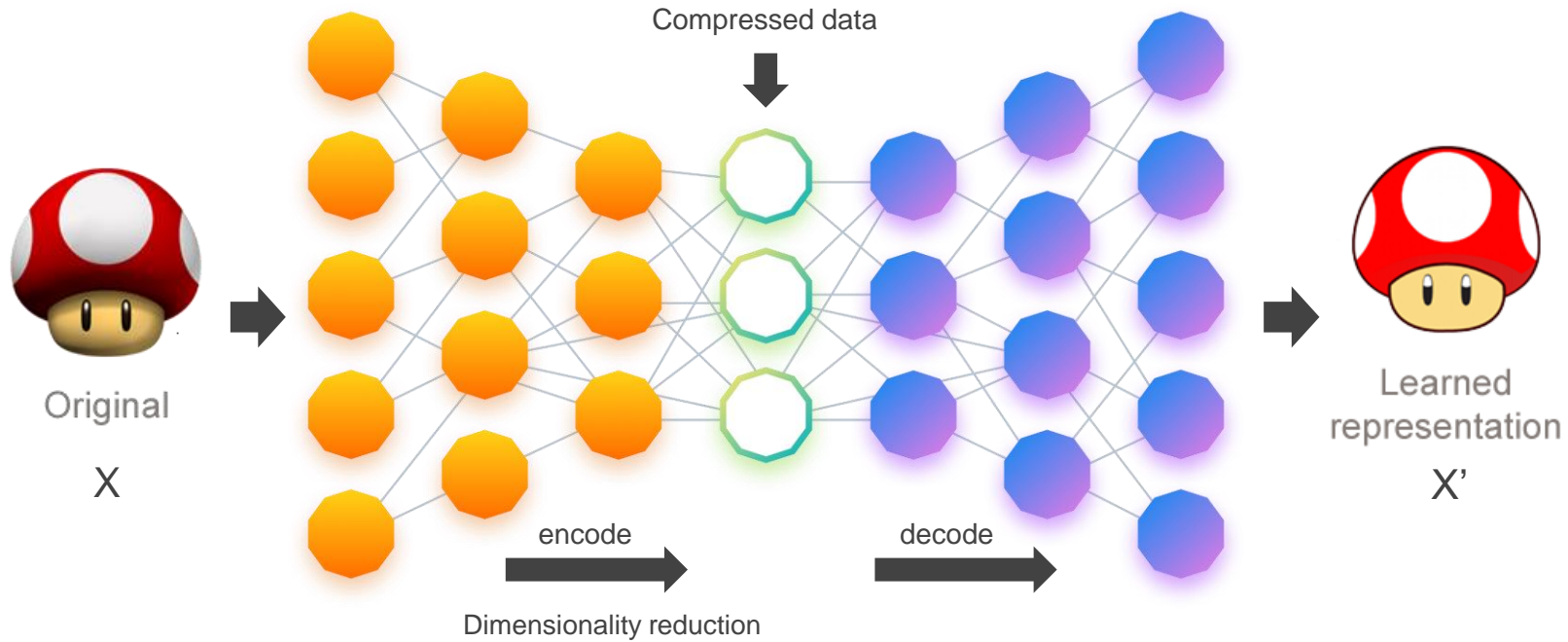


# Outlier Detection Range Numbers

Deep Learning - Autoencoders



# Autoencoders



Similar data preparation steps

- Range
- Range Level
- Signaling Events CV
- AB Ratio
- Rotation Score
- Interval consistency
- (...)

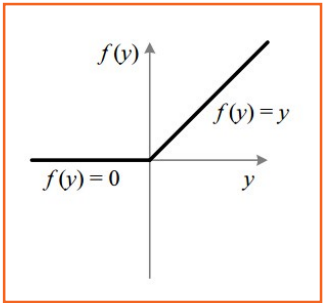
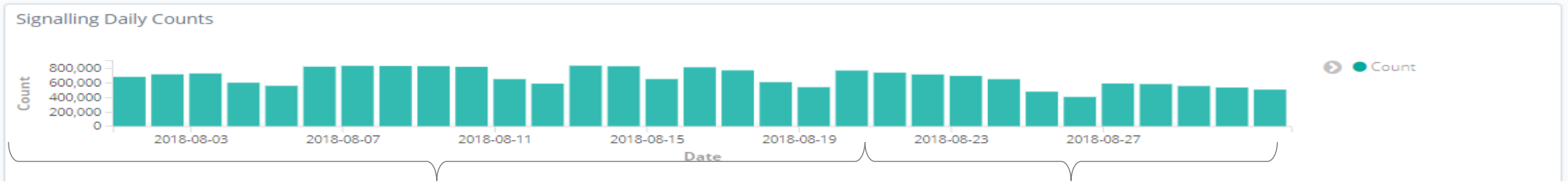


standardized scale ( $z$ )

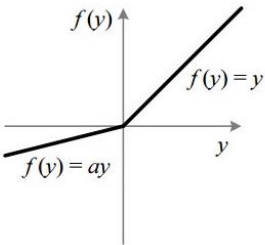
$$z = \frac{x - \mu}{\sigma}$$



# Autoencoders - Results



ReLU



Leaky ReLU

Fit

13-8-2-8-13

Optimizer	Adam
Loss	MSE
Act. Function	ReLU
Epochs	50

Predict



**MSE**



Tukey (k=3)  
83 459 cases

Q(0.99)  
6807 cases

20 / 21 known bypass (95%)  
2 / 52 known FP (2%)

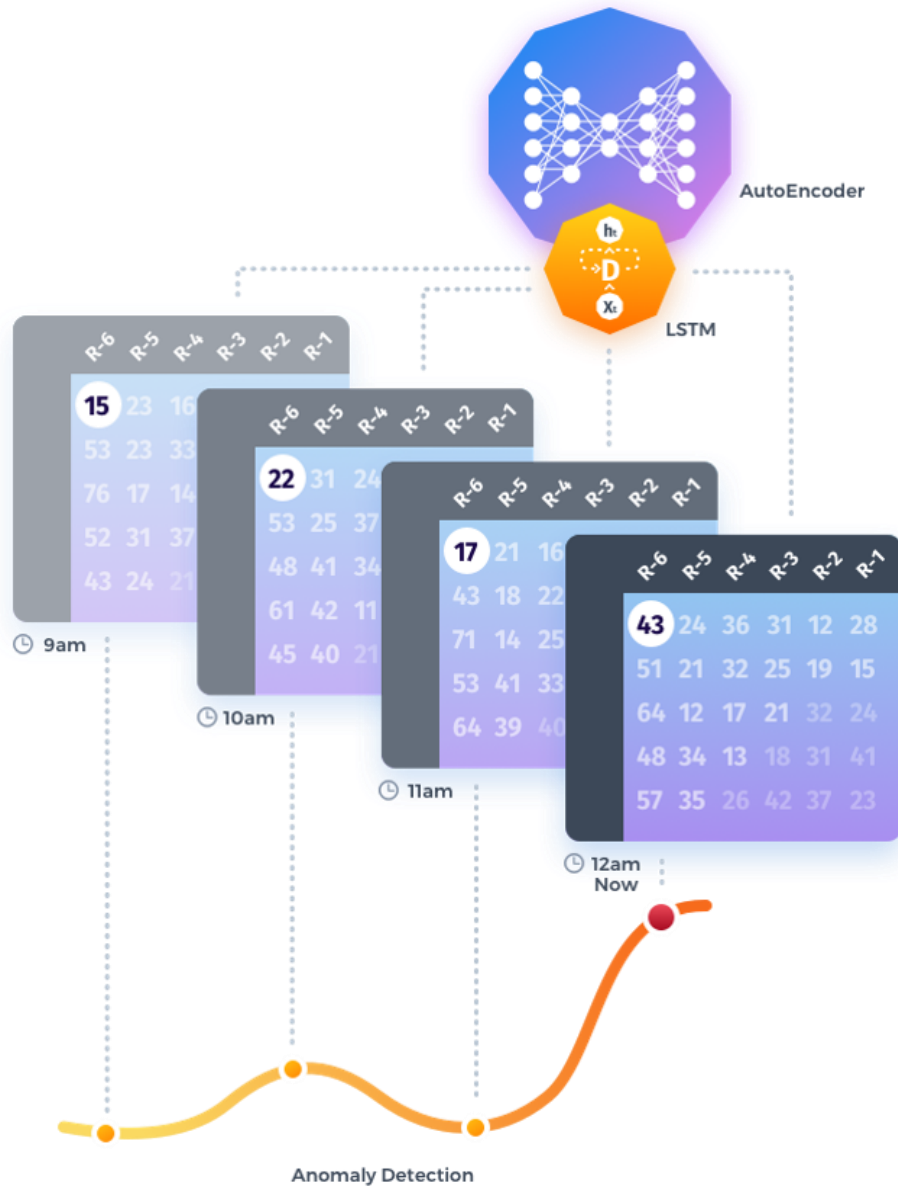
16 / 21 known bypass (76%)  
0 / 52 known FP (0%)

# Outlier Detection Range Numbers

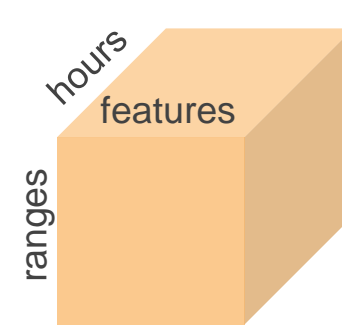
Deep Learning – LSTM + AE



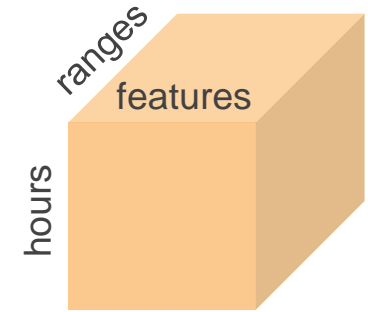
# Stacked LSTM + Autoencoders



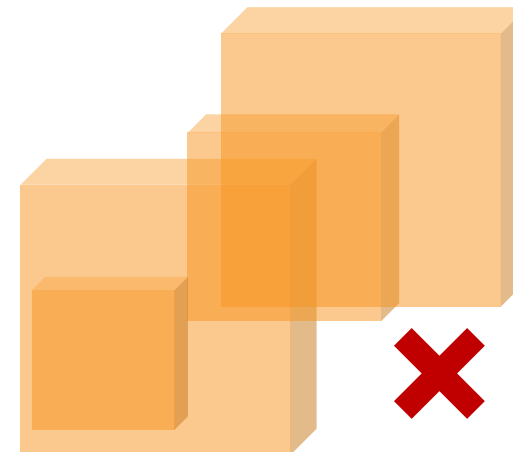
## 3D Approach



Sequence of hours  
(scaled)

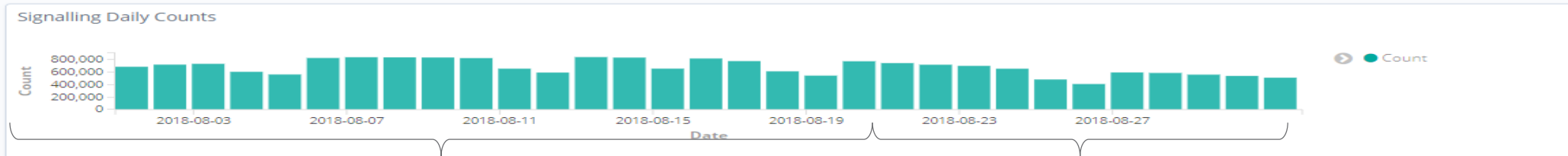


Sequence of ranges  
(scaled)



**✗** Different input shape

# Stacked LSTM + Autoencoders - Results

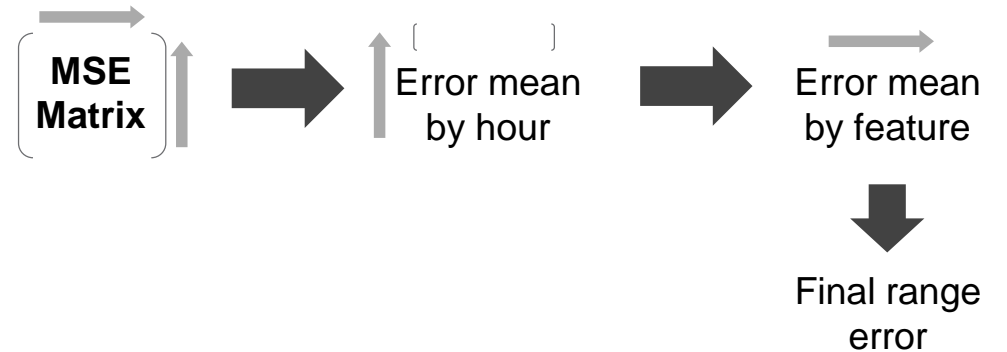


Fit

Example in Keras

```
inputs = Input(shape=(X.shape[1], X.shape[2]))
L1 = LSTM(16, activation='relu', return_sequences=True,
          kernel_regularizer=regularizers.l2(0.00))(inputs)
L2 = LSTM(4, activation='relu', return_sequences=False)(L1)
L3 = RepeatVector(X.shape[1])(L2)
L4 = LSTM(4, activation='relu', return_sequences=True)(L3)
L5 = LSTM(16, activation='relu', return_sequences=True)(L4)
output = TimeDistributed(Dense(X.shape[2]))(L5)
model = Model(inputs=inputs, outputs=output)
```

Predict



Optimizer	Adam
Loss	MSE
Act. Function	ReLU
Epochs	50

Tukey (k=3)  
112 713 cases



Q(0.99)  
10 235 cases

21 / 21 known bypass (100%)

20 / 21 known bypass (95%)

10 / 52 known FP (20%)

8 / 52 known FP (15%)



# Conclusions



# Conclusions

## Statistical with scoring approach

- Fast (18 secs for 10 days)
- Good results with **low cases**
- **Keeps the history** of daily studies to help end user to understand what is going on
- Simple and **easy to explain**
- Overfitting can occur

## Autoencoders

- Can learn new feature representation
- Relatively simple to understand
- Hyper parameter tuning
- Good results but only with a **high number of cases**
- Not so easy to explain results (feature importance)

## PCA

- Fast
- Good results but only with a **high number of cases**
- Not easy to explain results

## LSTM + Autoencoders

- Can learn new feature representation
- Sequence learning capability
- Good results require **very high number of cases (high FPs)**
- Very high computational cost (fit)
- Complex
- Hard to explain results

# THANK YOU

[raphael.espanha@wedotechnologies.com](mailto:raphael.espanha@wedotechnologies.com)

**Know** the unknown ...





# Q & A