



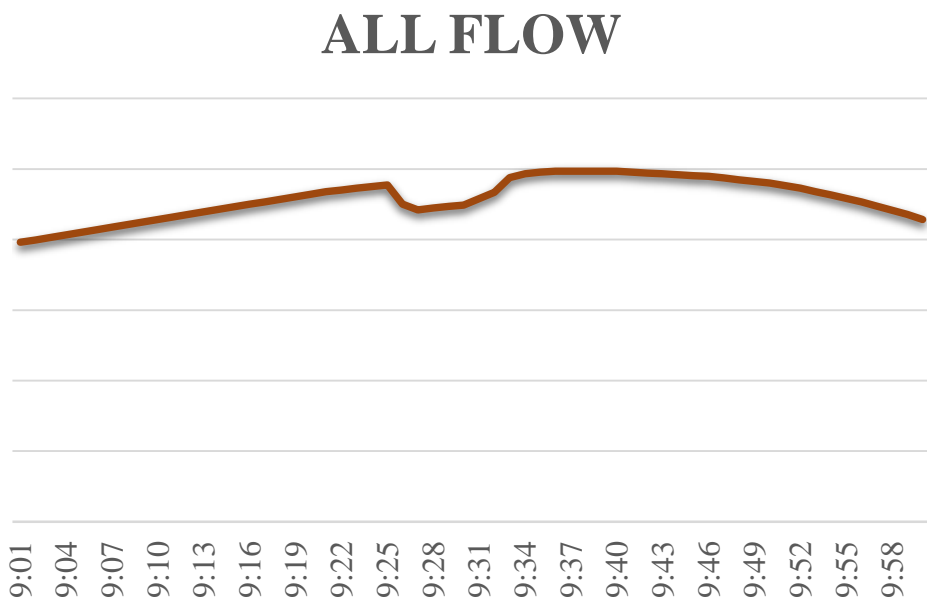
# Efficient trouble shooting of service failures with multi-tag data analysis

---

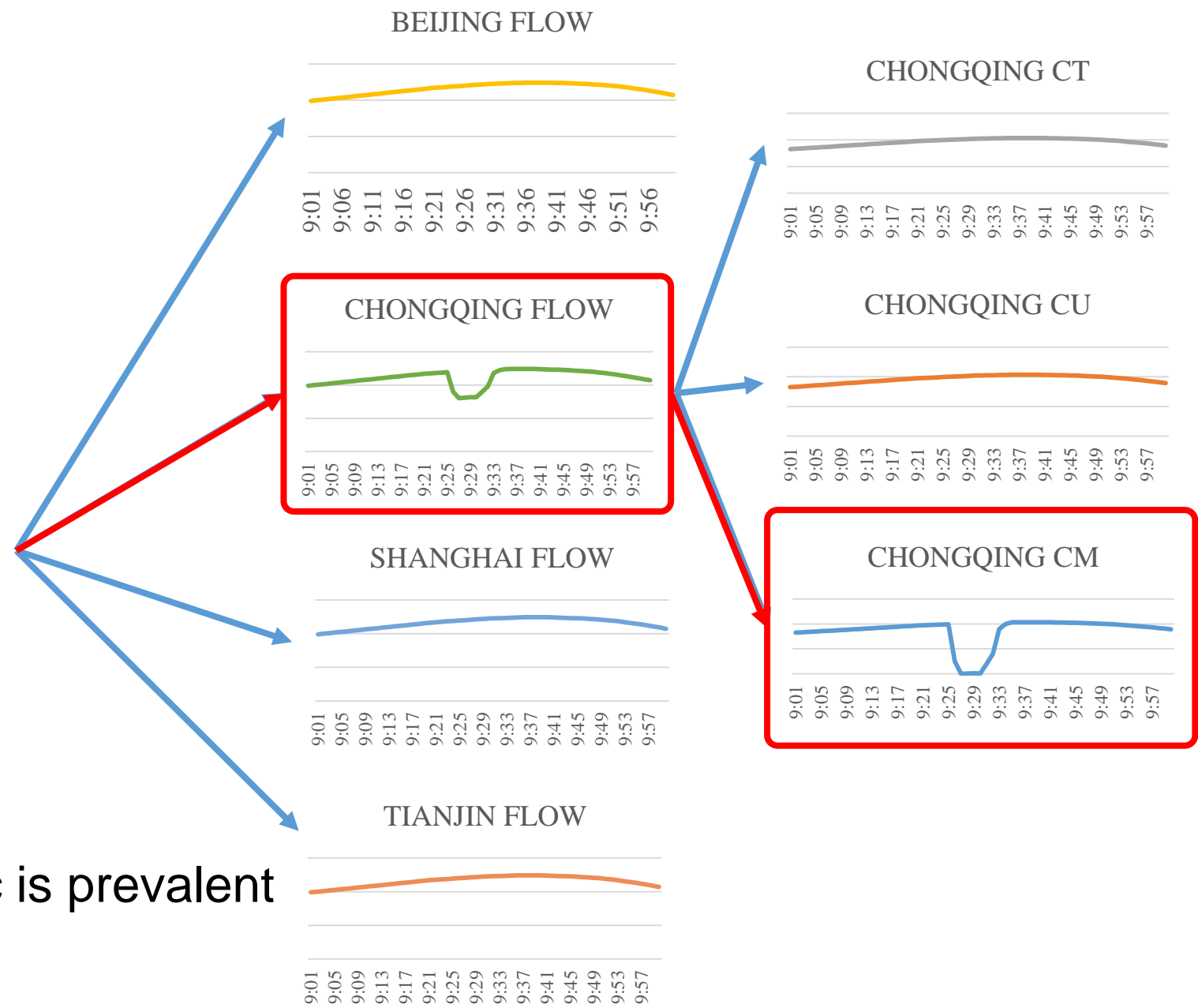
Cao, Xuan  
Baidu SRE

# What's the trouble?

- Facing to the problem

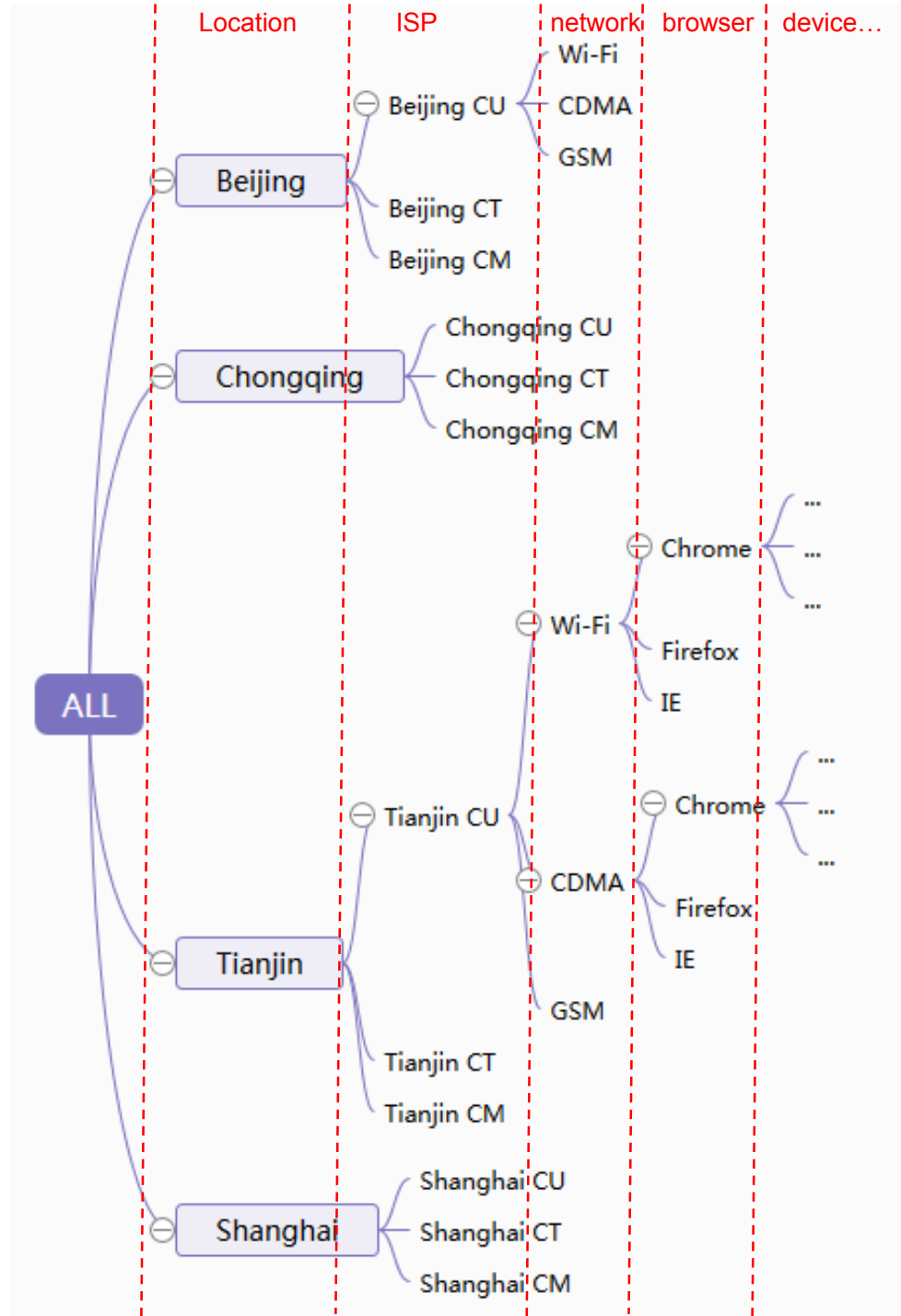


- Troubles affecting partial traffic is prevalent



# What's the trouble?

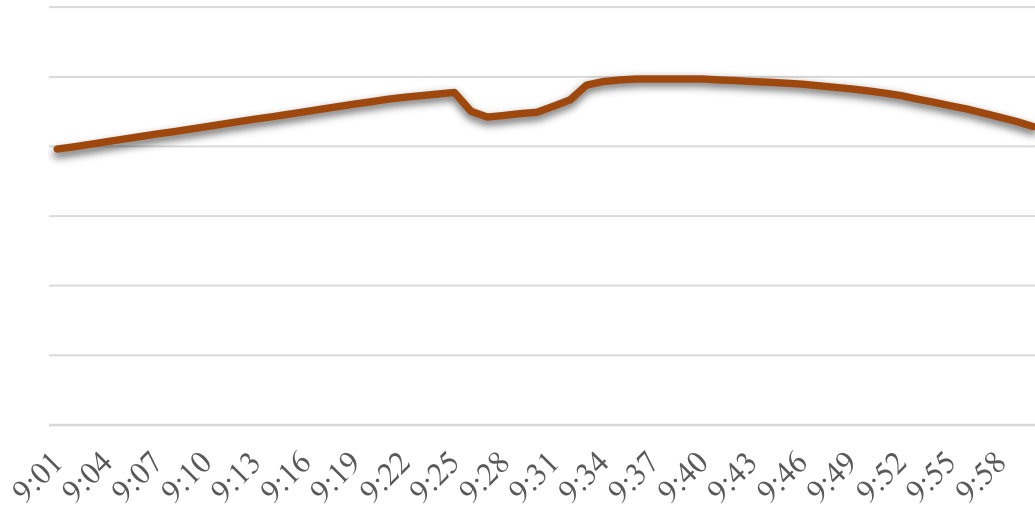
- Descartes accumulation set of all dimensions
- Lots of searching branches(waste of time)
- Need to narrow down the search scope
- Prune——Depending on SRE's experience



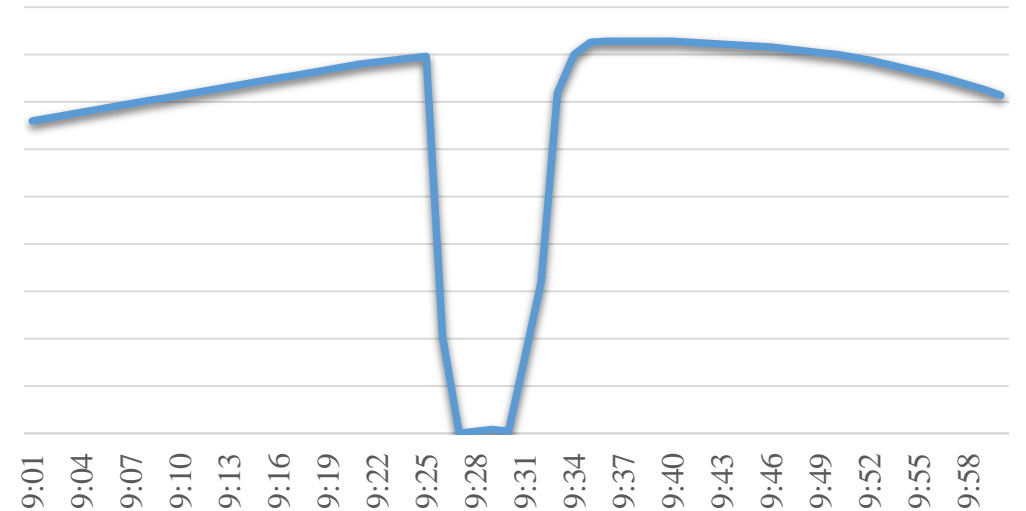
# Ideal result

- We directly got the answer

## ALL FLOW



## CHONGQING CM



# Our solution

- Pick a key indicator
- Procedure
  - Feature extraction —— assigning tags
  - Unsupervised anomaly detection
  - Entropy-based dimension reduction

# Assigning the tags

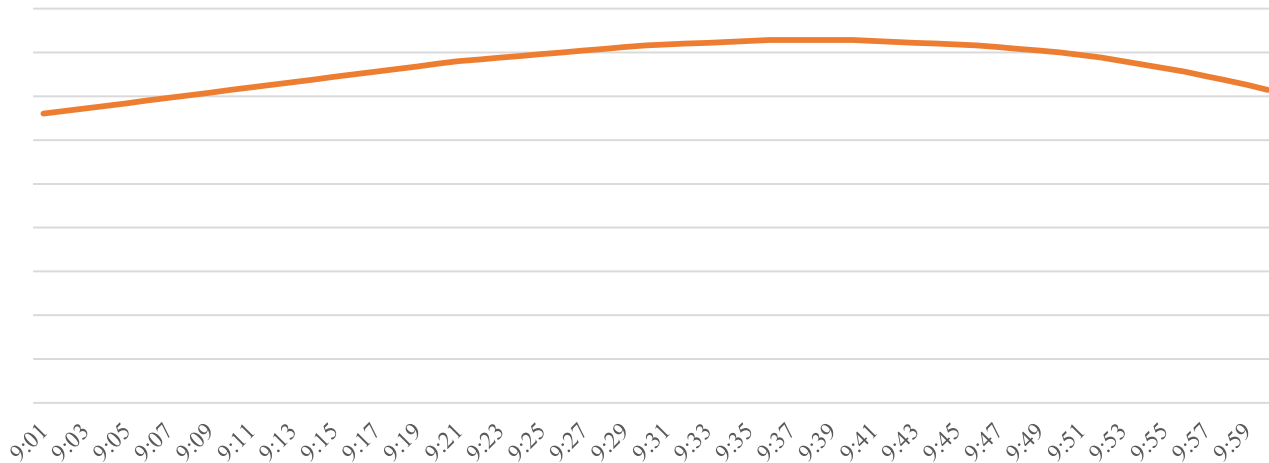
- Client-side tags
  - traffic location
  - browser type
  - access network standard
  - device type
  - ...
- Server-side tags
  - Service IDC
  - API Version
  - API type
  - ...

Query Word	Source area	ANS	browser	device	ISP	IDC	...
Driverless car	China	CDMA	Safari	Cell phone	CUCC	IDC-A	...
Sweater	Singapore	Wi-Fi	Chrome	Cell phone	Singtel	IDC-B	...
Machine learning	USA	Wi-Fi	Chrome	pad	T-Mobile	IDC-A	...
Forbidden city	China	Wi-Fi	Firefox	PC	CMCC	IDC-B	...
Pancake rolled with crisp fritter	Singapore	LTE	Safari	Cell phone	M1	IDC-A	...
...	...	...	...	...	...	...	...

# Assigning the tags

- traffic location ( country/province/city/etc... )
- browser type ( chrome/safari/firefox/etc... )
- access network standard ( Wi-Fi/CDMA/LTE/etc... )
- device type ( PC/laptop/pad/cell phone/etc... )
- ...

PV from China & ANS is Wi-Fi & device type is cell phone  
Time series trend diagram



# Unsupervised anomaly detection

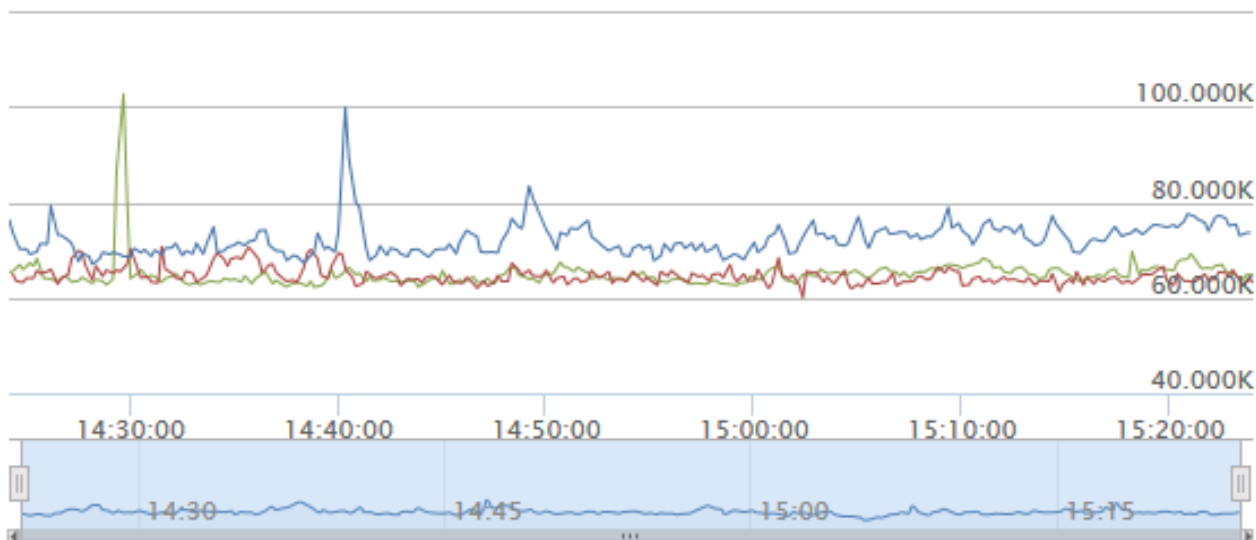
- Each algorithm works very well for certain types of indicators
- Unsupervised training to get thresholds for all finest dimension combinations
  - why unsupervised? thousands of combinations
  - train thresholds based on history data
  - use latency as an example



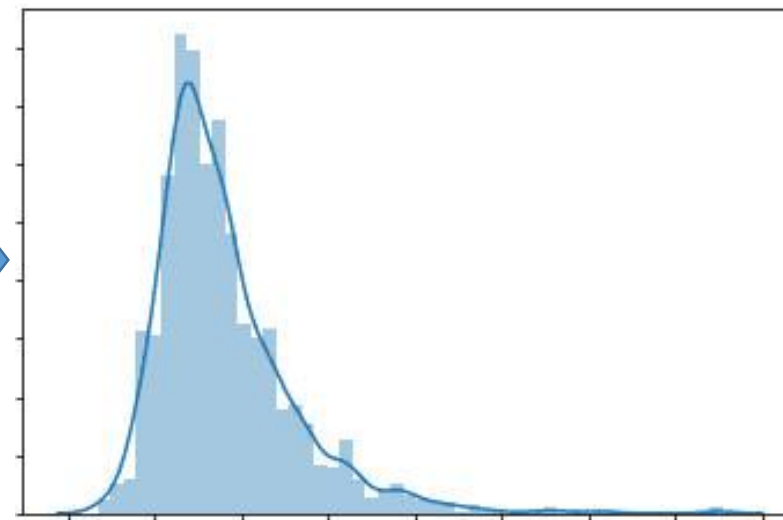
# Unsupervised anomaly detection

- e.g. Anomaly detection on service latency
  - KEY- how to determine an appropriate delay threshold
  - build a probability distribution for latency values
  - usually single-peak distribution on the histogram

Time series service latency

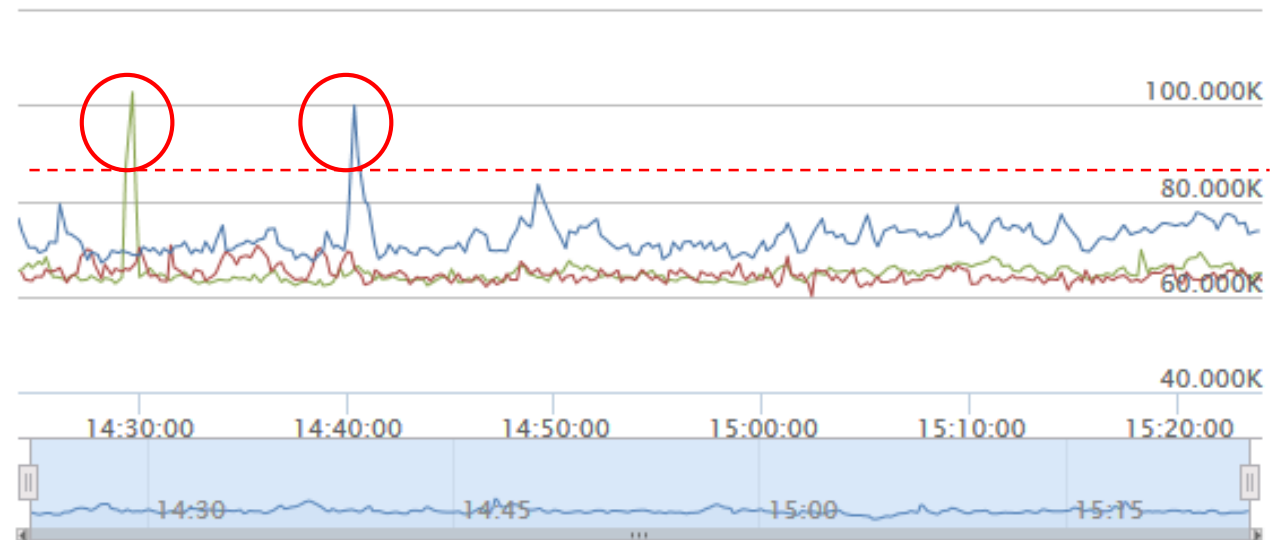
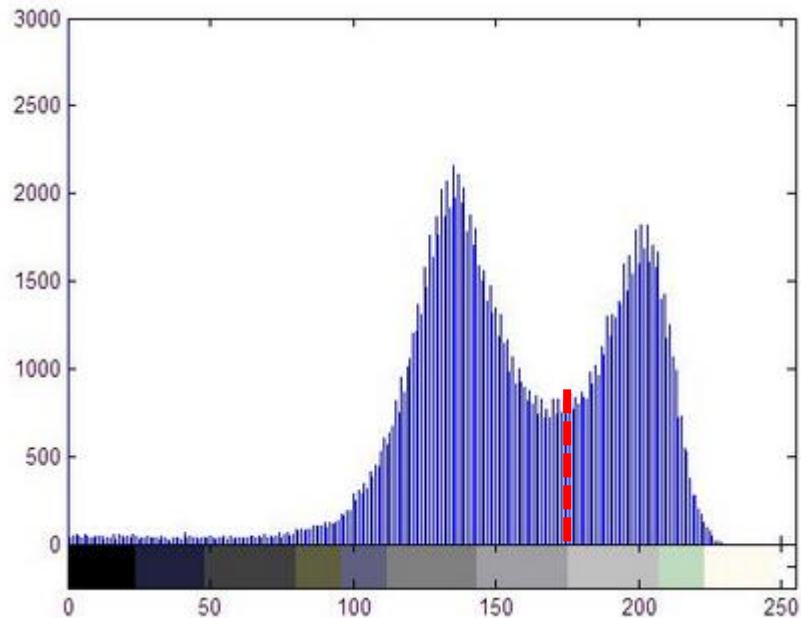


Histogram of service latency



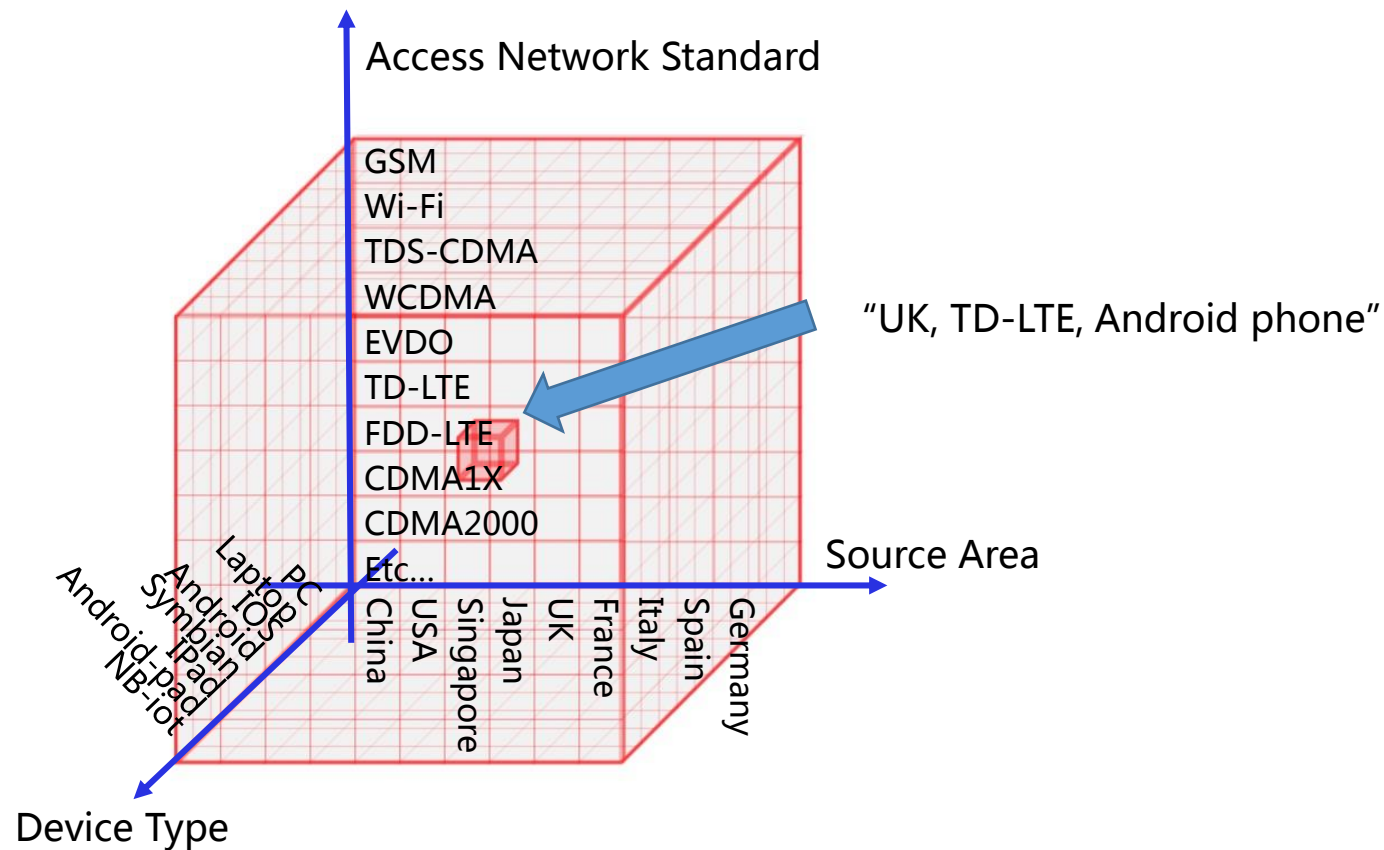
# Unsupervised anomaly detection

- Anomaly detection on service latency
  - Two/multi-peak distribution when failure happens
  - Maximize between-class scatter -> Threshold



# Entropy-based dimension reduction

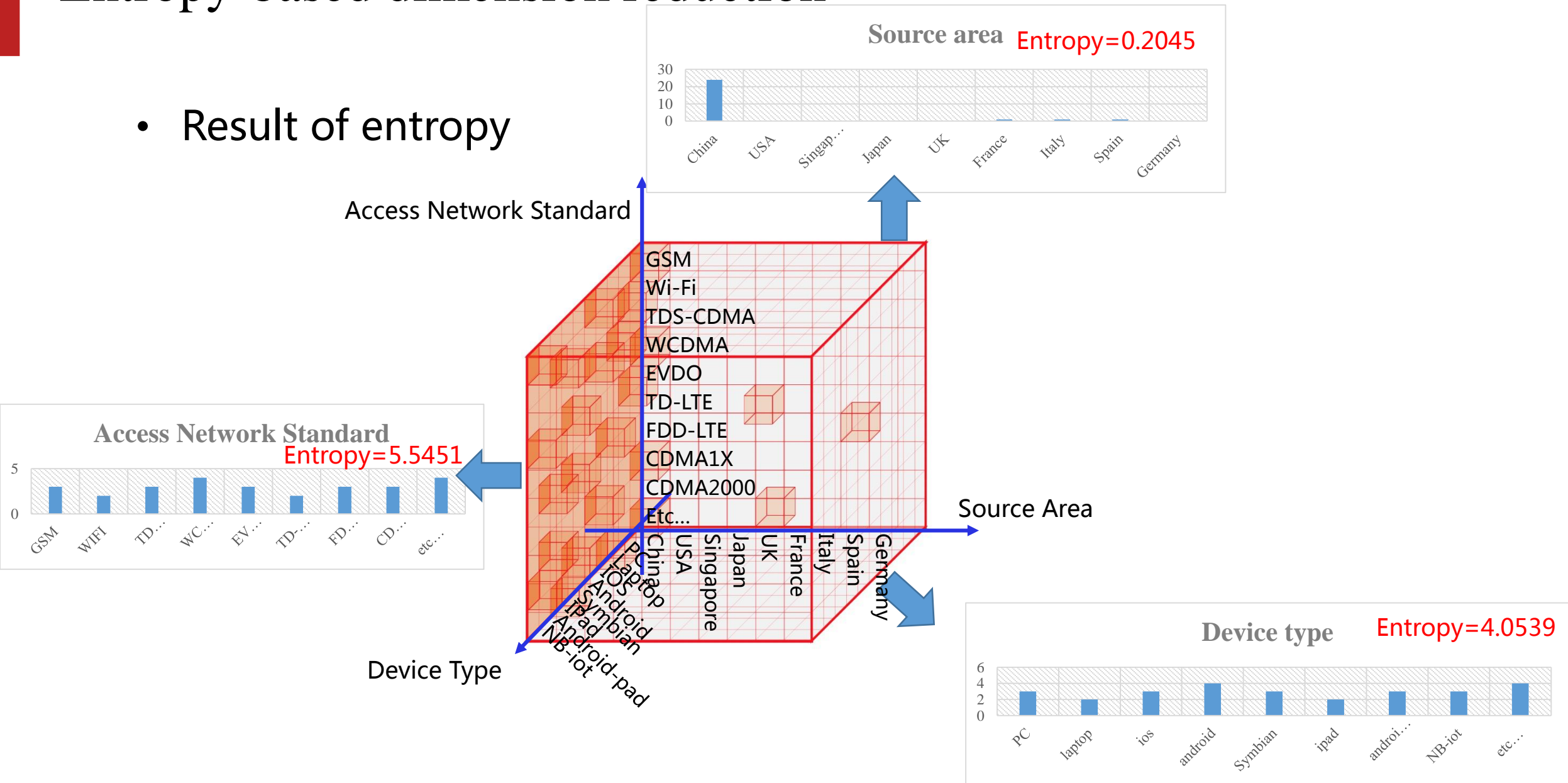
- Error cube





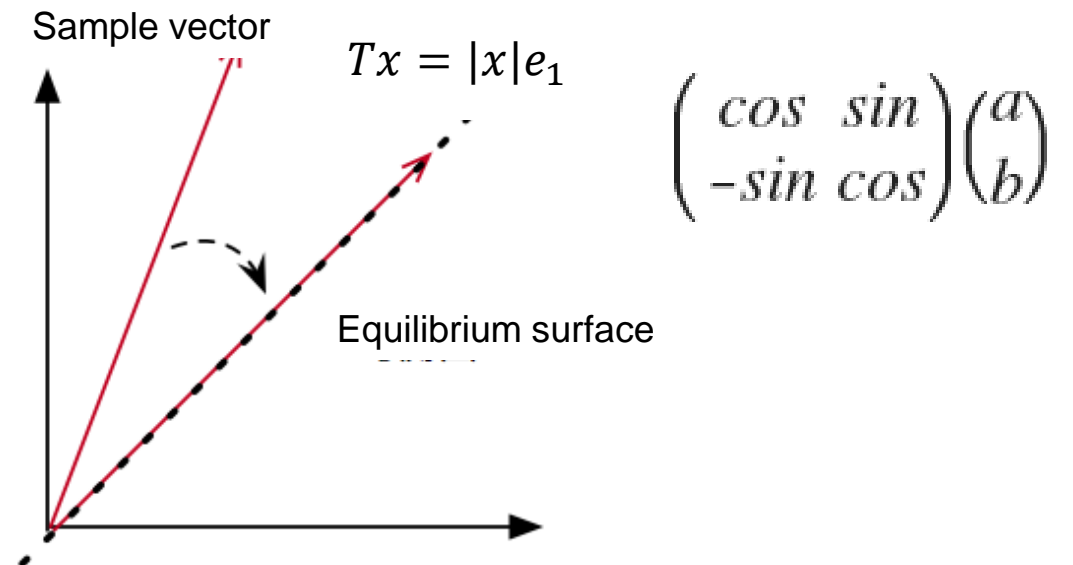
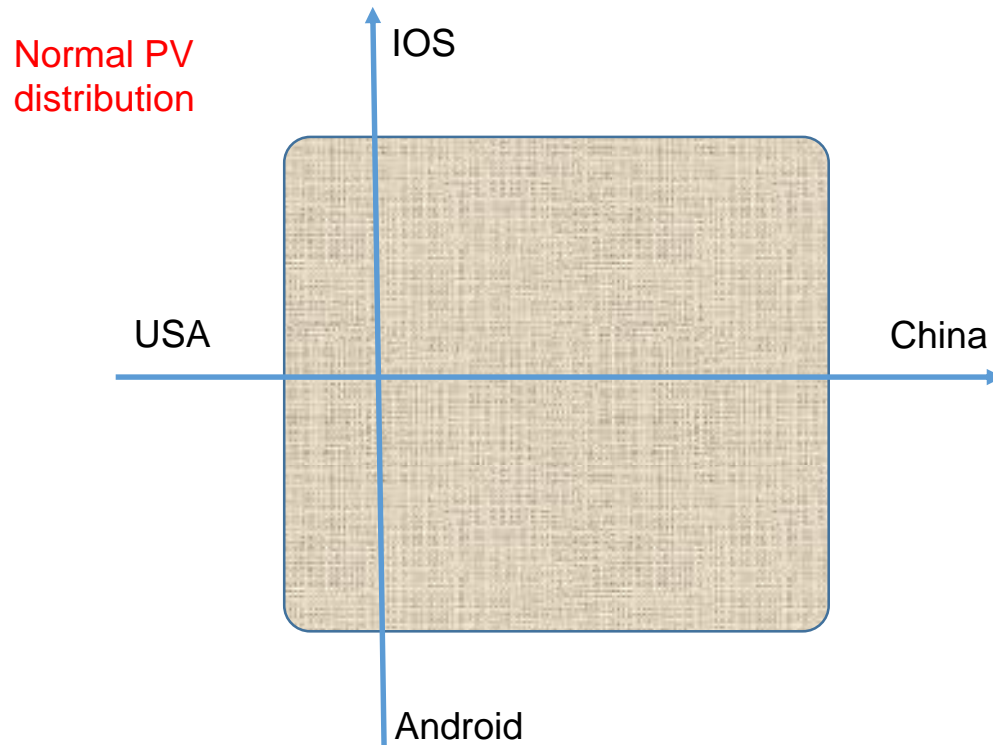
# Entropy-based dimension reduction

- Result of entropy



# Unbalanced base distribution

- Mostly for metrics such as PV/PVLOST
- Gives transformation to convert base distribution to uniform



- transform error distribution in the same way



# Summary

- What's the trouble?
  - feature extraction
  - unsupervised anomaly detection
  - entropy-based dimension reduction
- 
- Q & A
  - [caoxuan@baidu.com](mailto:caoxuan@baidu.com)





THANK YOU!

