

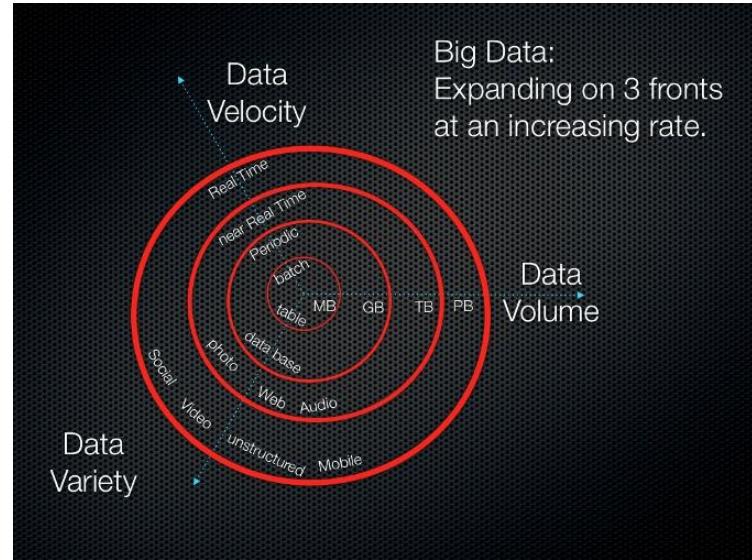
Big data and Data Science

Chantana Chantrapornchai
Dept. of Computer Engineering
Kasetsart University

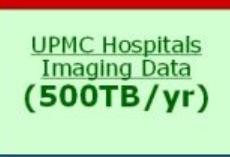


What is Big Data? - 3 V's

- Big Volume
 - With simple (SQL) analytics
 - With complex (non-SQL) analytics
- Big Velocity
 - How fast data is processed
- Big Variety
 - Large number of diverse data sources to integrate



<http://whatis.techtarget.com/definition/3Vs>

 Particle Physics Large Hadron Collider (15PB)	 Human Genomics (7000PB) 1GB / person 200PB+ captured 200% CAGR	 http://www.intel.com World Wide Web (~1PB) http://www.intel.com	 Wikipedia (10GB) 100% CAGR http://www.wikipedia.org
 Annual Email Traffic, no spam (300PB+)	 INTERNET ARCHIVE (1PB+)	 Estimated On-line RAM in Google (8PB)	 Personal Digital Photos (1000PB+) 100% CAGR
 200 of London's Traffic Cams (8TB/day)	 2004 Walmart Transaction DB (500TB)	 Typical Oil Company (350TB+)	 Merck Bio Research DB (1.5TB/qtr)
 UPMC Hospitals Imaging Data (500TB/yr)	 MIT Babyltalk Speech Experiment (1.4PB)	 Terashake Earthquake Model of LA Basin (1PB)	 One Day of Instant Messaging in 2002 (750GB)
Total digital data to be created this year 270,000PB (IDC)			

Compound Annual Growth Rate: CAGR
<http://www.investopedia.com/terms/c/cagr.asp>

12 big facts about big data

1. If you stacked a pile of CD-ROMs on top of one another until you'd reached the current global **storage capacity** for digital information – about 295 exabytes – it would stretch 80,000 km beyond the moon. **จำนวน storage ที่ใช้**

2 Every hour, enough information is consumed by **internet traffic** to fill 7 million DVDs. Side by side, they'd scale Mount Everest 95 times. **จำนวนข้อมูลที่ไหลบน Internet**

3 247 billion e-mail messages are sent each day... up to 80% of them are spam.

4 By 2020, IT departments will be looking after 10 x **more servers**, 50 x more data and 75 x more files. Meanwhile, the number of IT administrators keeping track of all that data growth will increase by 1.5 times. **จำนวน server ที่ใช้**

5 We can expect a 40-60 per cent projected **annual growth in the volume** of data generated, while media intensive sectors, including financial services, will see year on year data growth rates of over 120 per cent.

6 The world's 500,000+ **data centres** are large enough to fill 5,955 football fields. **ปริมาณที่เพิ่มขึ้นทุกปี**

7 775% of digital information is generated by individuals, whilst enterprises have liability for 80% of digital data at some point in its life. **ขนาด data center**

8 There are nearly as **many bits of information** in the digital universe as there are stars in our actual universe.

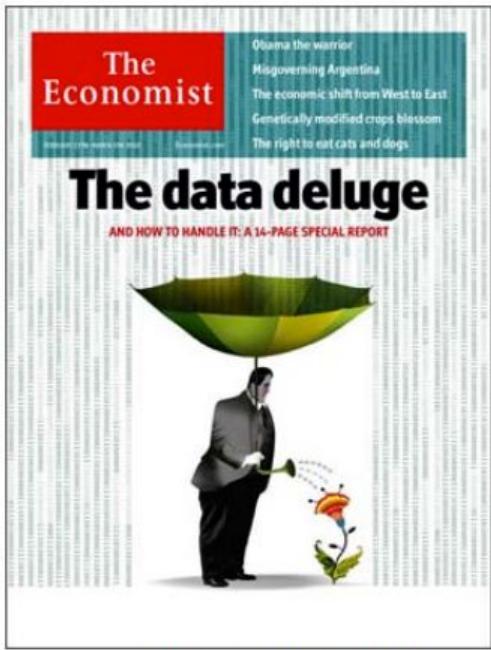
9 **Investment in digital enterprises** has increased 50% since 2005. **การลงทุนในหน่วยงานเพิ่มขึ้น 50%**

10 There are 30 billion **pieces of content shared on Facebook** every day. **จำนวนข้อมูลที่ shared บน facebook ทุกวัน**

11 In 2010, 28% of the digital universe required some level of security... not all of it had the level of security it required....

12 People wishing each other Happy New Year drove a 500% surge in smartphone data within just one year, according to 3UK whose customers used a whopping 80 terabytes (TB) on the 31st December 2011, compared to just 14 TBs on the same day in 2010. **ความต้องการ security**

จำนวน ข้อมูลที่ส่งบน mobile ในวันปีใหม่



The Economist, Feb 25, 2010

IN 2010 THE DIGITAL UNIVERSE WAS
1.2 ZETTABYTES

IN A DECADE THE DIGITAL UNIVERSE WILL BE
35 ZETTABYTES

90% OF THE DIGITAL UNIVERSE IS
UNSTRUCTURED

IN 2011 THE DIGITAL UNIVERSE IS
300 QUADRILLION FILES

WIRED **The New York Times** **Bloomberg Businessweek** **Forbes** **WALL STREET JOURNAL**

http://www.snia.org/sites/default/files2/ABDS2012/Tutorials/RobPeglar_Introduction_Analytics%20_Big%20Data_Hadoop.pdf

How big?

- $1 \text{ ZB} = 1000^7 \text{ bytes} = 10^{21} \text{ bytes} = 100000000000000000000000 \text{ bytes} = 1000 \text{ exabytes} = 1 \text{ billion tera bytes.}$

Multiples of bytes			V*T*E
Decimal		Binary	
Value	Metric	Value	JEDEC
1000	kB kilobyte	1024	KB kilobyte
1000^2	MB megabyte	1024^2	MB megabyte
1000^3	GB gigabyte	1024^3	GiB gibibyte
1000^4	TB terabyte	1024^4	TiB tebibyte
1000^5	PB petabyte	1024^5	PiB pebibyte
1000^6	EB exabyte	1024^6	EiB exbibyte
1000^7	ZB zettabyte	1024^7	ZiB zebibyte
1000^8	YB yottabyte	1024^8	YiB yobibyte

Orders of magnitude of data

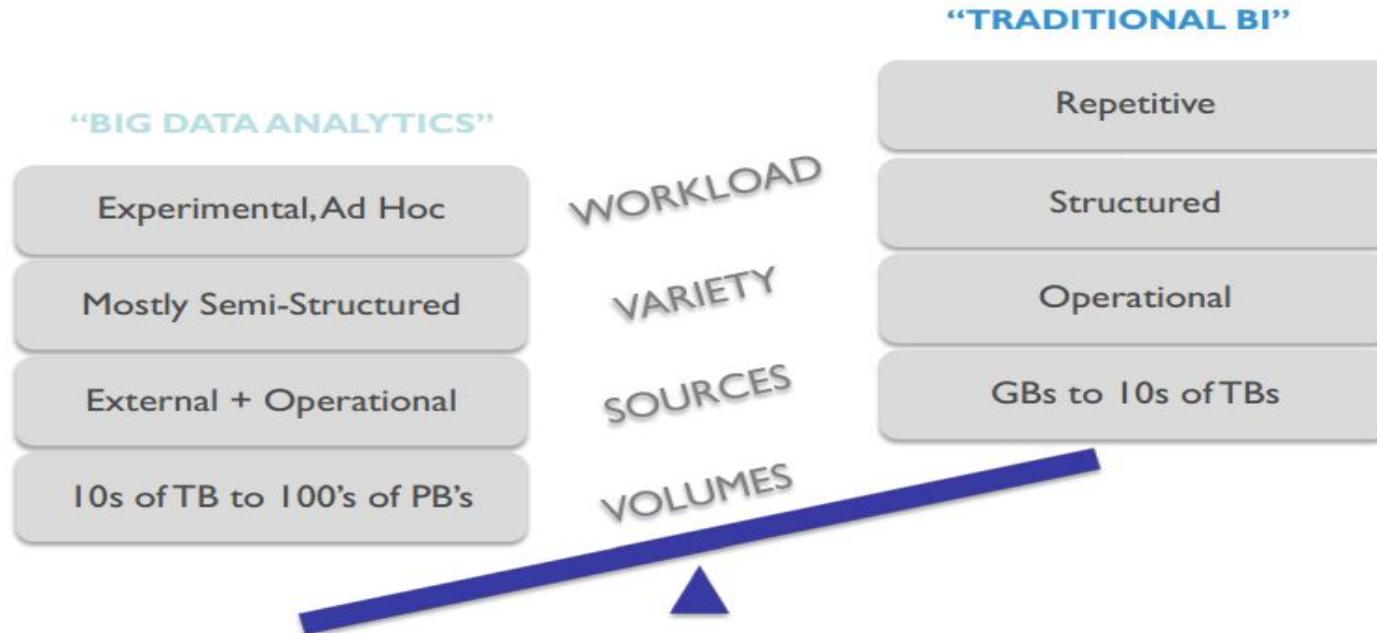
<http://en.wikipedia.org/wiki/Zettabyte>



And Future Growth is Staggering



Processing big data: Differences from traditional BI



Ten commons applications

1. Modeling true risk
2. Customer churn analysis
3. Recommendation engine
4. Ad targeting
5. PoS transaction analysis
6. Analyzing network data to predict failure
7. Threat analysis
8. Trade surveillance
9. Search quality
10. Data “sandbox”



Retail

- CRM – Customer Scoring
- Store Siting and Layout
- Fraud Detection / Prevention
- Supply Chain Optimization



Advertising & Public Relations

- Demand Signaling
- Ad Targeting
- Sentiment Analysis
- Customer Acquisition



Financial Services

- Algorithmic Trading
- Risk Analysis
- Fraud Detection
- Portfolio Analysis



Media & Telecommunications

- Network Optimization
- Customer Scoring
- Churn Prevention
- Fraud Prevention



Manufacturing

- Product Research
- Engineering Analytics
- Process & Quality Analysis
- Distribution Optimization



Energy

- Smart Grid
- Exploration



Government

- Market Governance
- Counter-Terrorism
- Econometrics
- Health Informatics



Healthcare & Life Sciences

- Pharmaco-Genomics
- Bio-Informatics
- Pharmaceutical Research
- Clinical Outcomes Research

Categories of Big Volumes

- With no analytics
- With analytics

Big volumes – no analytics

- Well addressed by **data warehouse** crowd
- Who are pretty good at SQL analytics on
 - Hundreds of nodes
 - Petabytes of data
- Column stores may be a good solution.
 - Factor of **50 or so** faster than row stores

Big data- big analytics

- Complex math operations (machine learning, clustering, trend detection,)
 - the world of the “quants” (Quantitative analyst)
 - Mostly specified as linear algebra on array data
- A dozen or so common ‘inner loops’
 - Matrix multiply
 - QR decomposition
 - SVD decomposition (Singular value decomposition)
 - Linear regression

Big Velocity

การซื้อขายใน
ตลาดหุ้น

- Trading volumes going through the roof on Wall Street
 - breaking infrastructure
- Sensor tagging of {cars, people, ...} creates a firehose to ingest
- The web empowers end users to submit transactions
 - sending volume through the roof
- PDAs lets them submit transactions from anywhere....

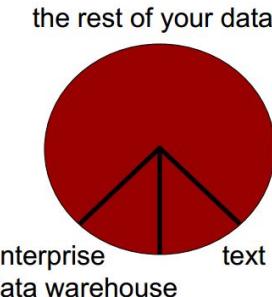
Big Variety

Typical enterprise has 5,000 operational systems

- Only a few get into the data warehouse
- What about the rest?

And what about all the rest of your data?

- Spreadsheets
- Access databases
- Web pages
- And public data from the web?



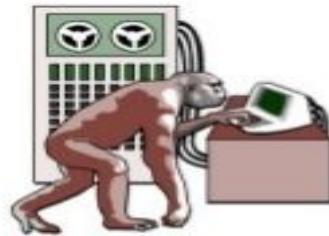
All data integrations

Information History

*Business Intelligence.
Great variety of visual resources
to analyze data*

*Beginning of the use of DBs
and basic reports*

Data was not stored



Data analysis profits:

- Competitive advantages
- Customer satisfaction evaluation
- Business process improvement
- Increase sales
- ...

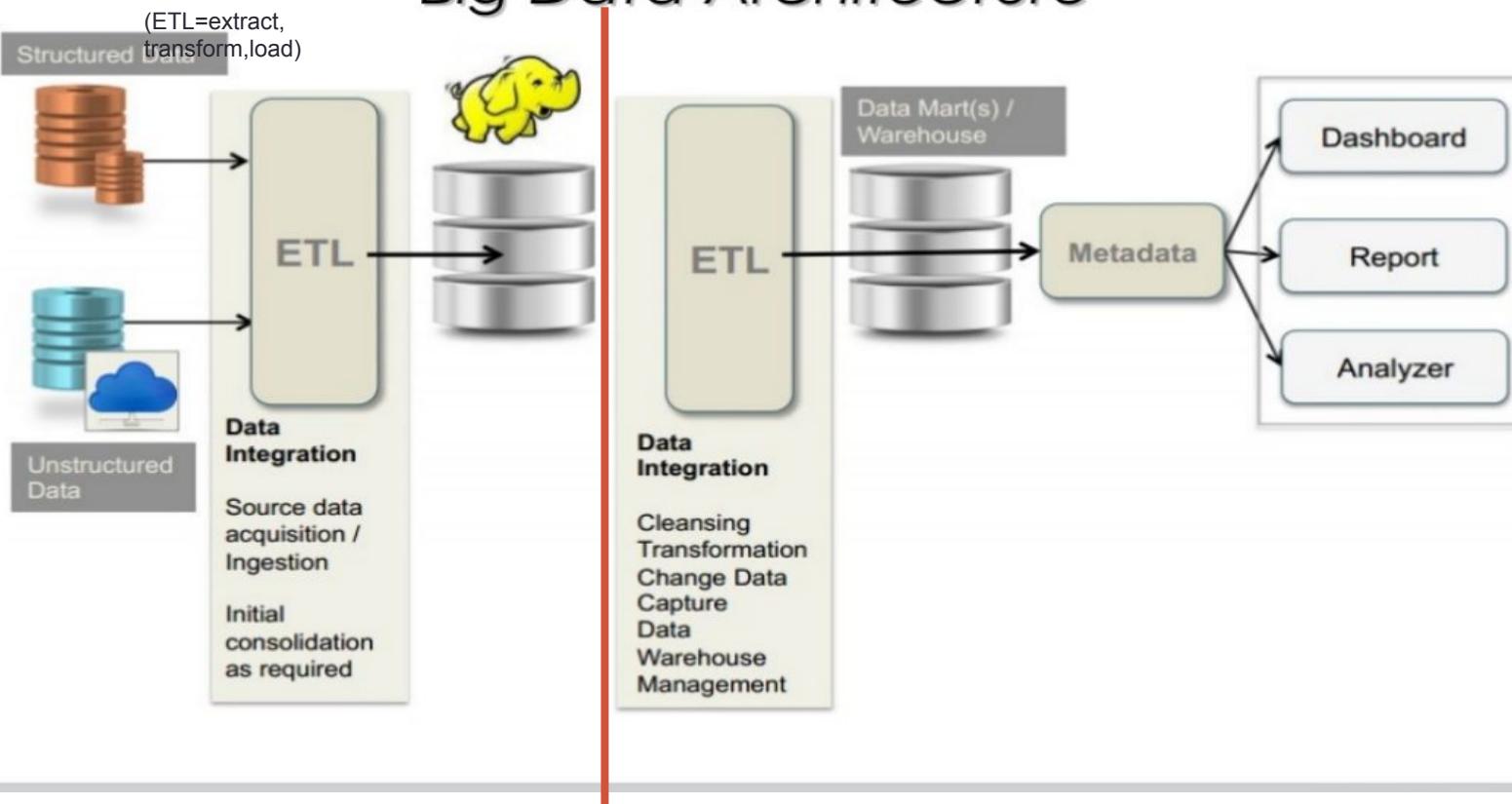


New data analysis techniques and processes

- New BI solutions
- New visual resources
- New data sources
- Cloud solutions
- Latest trends
 - Social Intelligence
 - Mailing intelligence
- ...

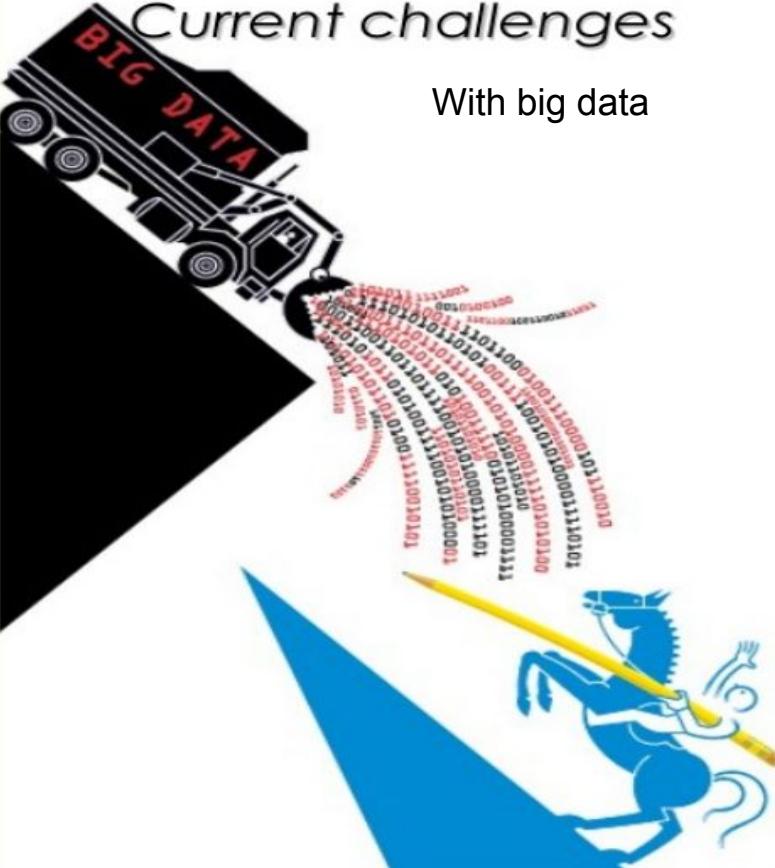


Big Data Architecture



Current challenges

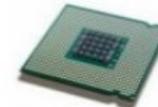
With big data



■ Scalability

■ Vertical

- + CPU
- + RAM



■ Horizontal

- More nodes



■ Data types

- Structured
- Unstructured

Data types

■ Structured

A *data structure* is a particular way of storing and organizing data in a computer so that it can be used efficiently.

List: http://en.wikipedia.org/wiki/List_of_data_structures

Primitive data types: Boolean, chart, float, double ...

■ Unstructured

Unstructured information refers to information that either does not have a pre-defined data model or is not organized in a pre-defined manner.

semi-structure....

Types of Big Data DBs. Not Only SQL (NoSQL)

- In response to these problems a NoSQL paradigm appeared.
- NoSQL **is not a substitute** for relational databases
 - Instead it is used in other **specific scenarios**
- Not all problems can be solved using a RDBMS
- Developer has a **range of possibilities** and can select the best to deal with a specific problem
- There are several NoSQL systems focusing on typical issues (scaling, increasing performance...) in a different way

EXAMPLE OF NOSQL

Document-oriented
databases



Columnar
databases

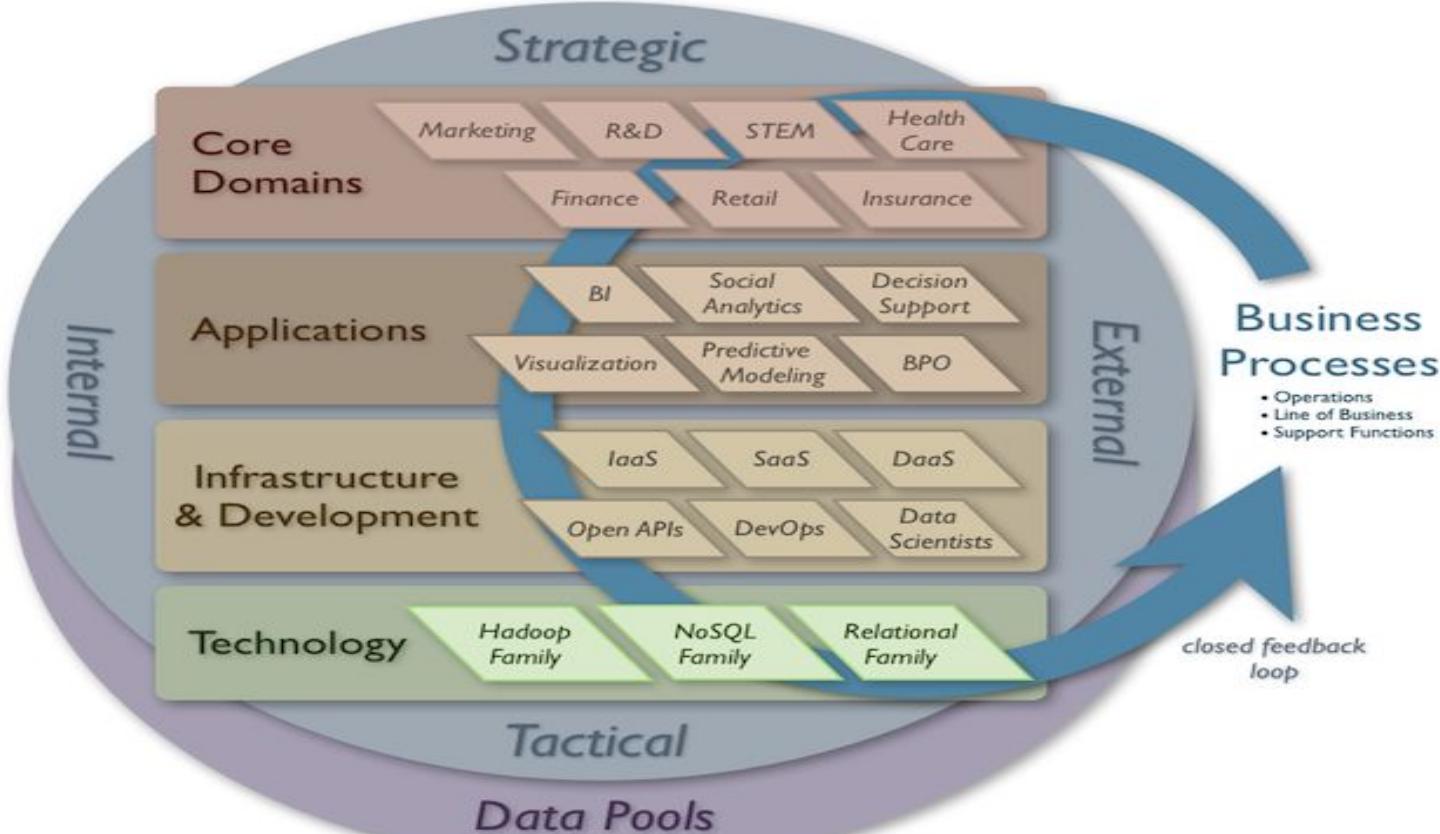


Object oriented
databases

Graph databases

Do not replace relational model. Specific scenarios.

Big Data in the Enterprise





Big Data: The Moving Parts

Increasing
Age & Maturity

Hadoop

Vertica

MapReduce

Esper

kdb

Greenplum

ETL

Netezza

ECL

Teradata

Fast Data

Hive

SciPy

Mahout

MATLAB

Revolution R

SPSS

AMPL

SAS

Big Analytics

unsupervised learning

social media analytics

sentiment analysis

predictive modeling

BPO

BI

network analysis

visualization

simulation

Deep Insight

Target goal

mass customization of services

quicker response to market trends

identifying real-time cost optimizations

faster, more accurate decision making

better and more holistic R&D

autonomic supply chain management

Business Objectives

From <http://blogs.zdnet.com/Hinchcliffe>

the growth of data will be exponential for the foreseeable future

terabytes

petabytes

exabytes

zettabytes

the amount of data stored by the average company today

Data Science

What is data science?

- Data Science is a field that comprises of everything that related to data cleansing, preparation, and analysis.
- Dealing with unstructured and structured data

Big Data: Big Data is humongous volumes of data that cannot be processed effectively with the traditional applications that exist.

Data Analytics: The science of examining raw data with the purpose of drawing conclusions.

It involves applying an algorithmic or mechanical process to derive insights.

Applications of data science

1. Internet search
2. Digital advertisement
3. Recommendation

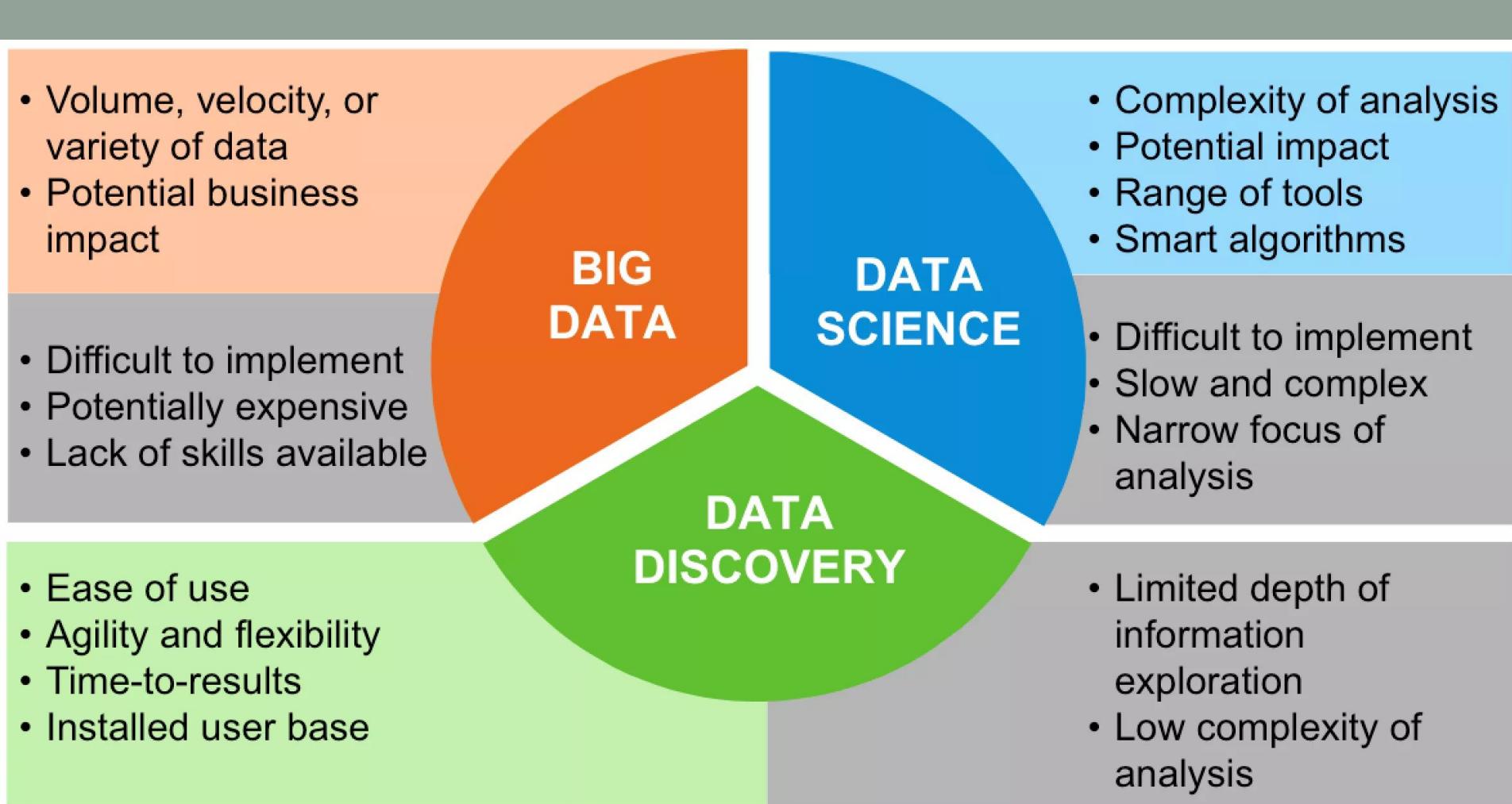
Application of Big data

Big data in communication: Gaining new subscribers, retaining customers, and expanding within current subscriber bases

Big Data for financial services: Credit card companies, retail banks, private wealth management advisories, insurance firms, venture funds, and institutional investment banks use big data for their financial services.

Applications of Data Analysis:

- Healthcare
- Tourism
- Gaming
- Energy management



Data Scientist

asking the
right questions

statistics

find/import/transform data
presentation skills
portfolio project

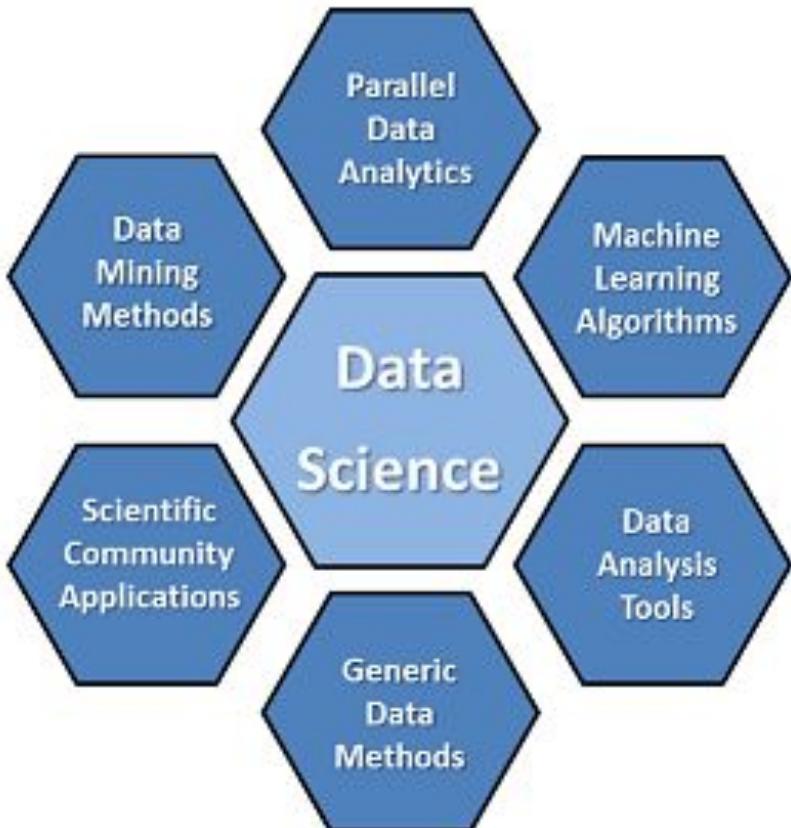
machine
learning

programming

Big-data Engineer

building a
product

streaming



Our case study: IoT data

Monitoring inefficiency of energy consumption in data center.

Chantana Chantrapornchai

HPCNC Laboratory, Dept. of Computer Engineering

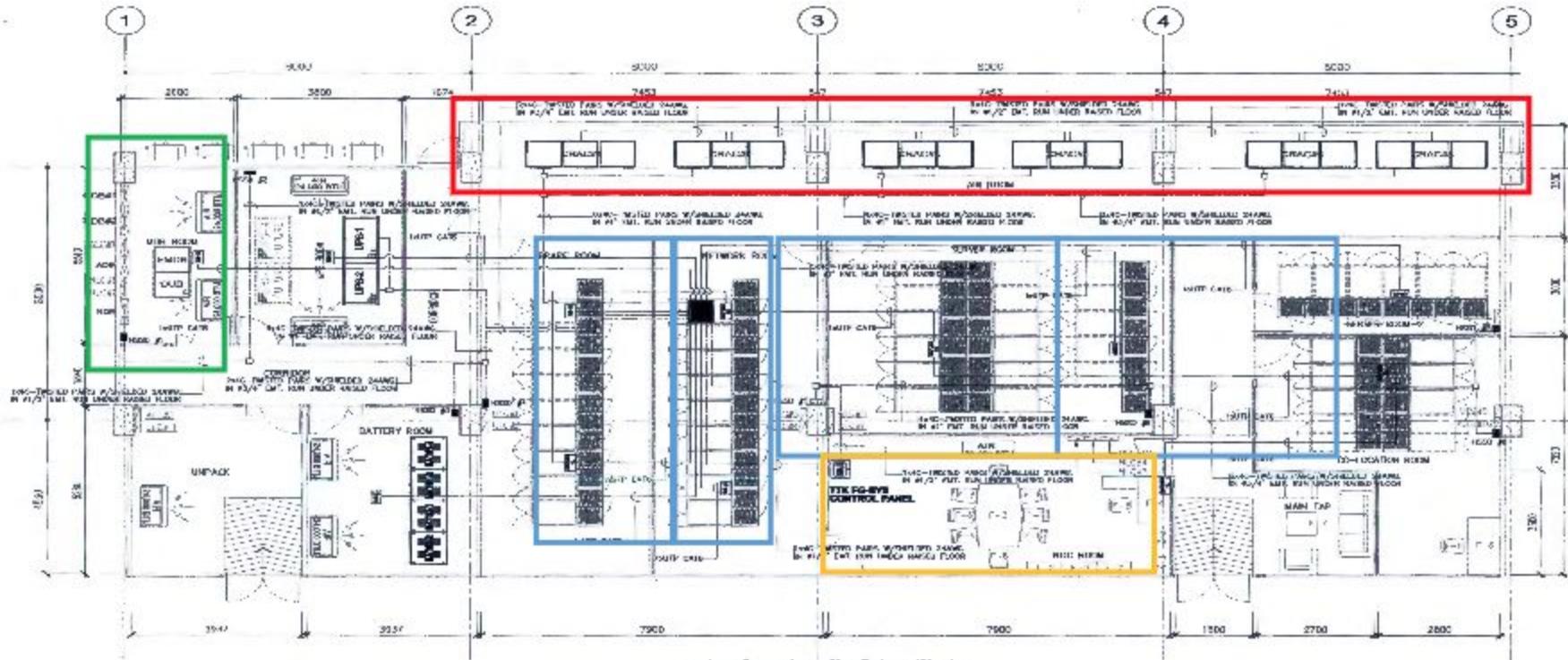
Kasetsart University

Withit Chatlatanakulchai, et.al.

CRVLAB , Dept. of Mechanical Engineering

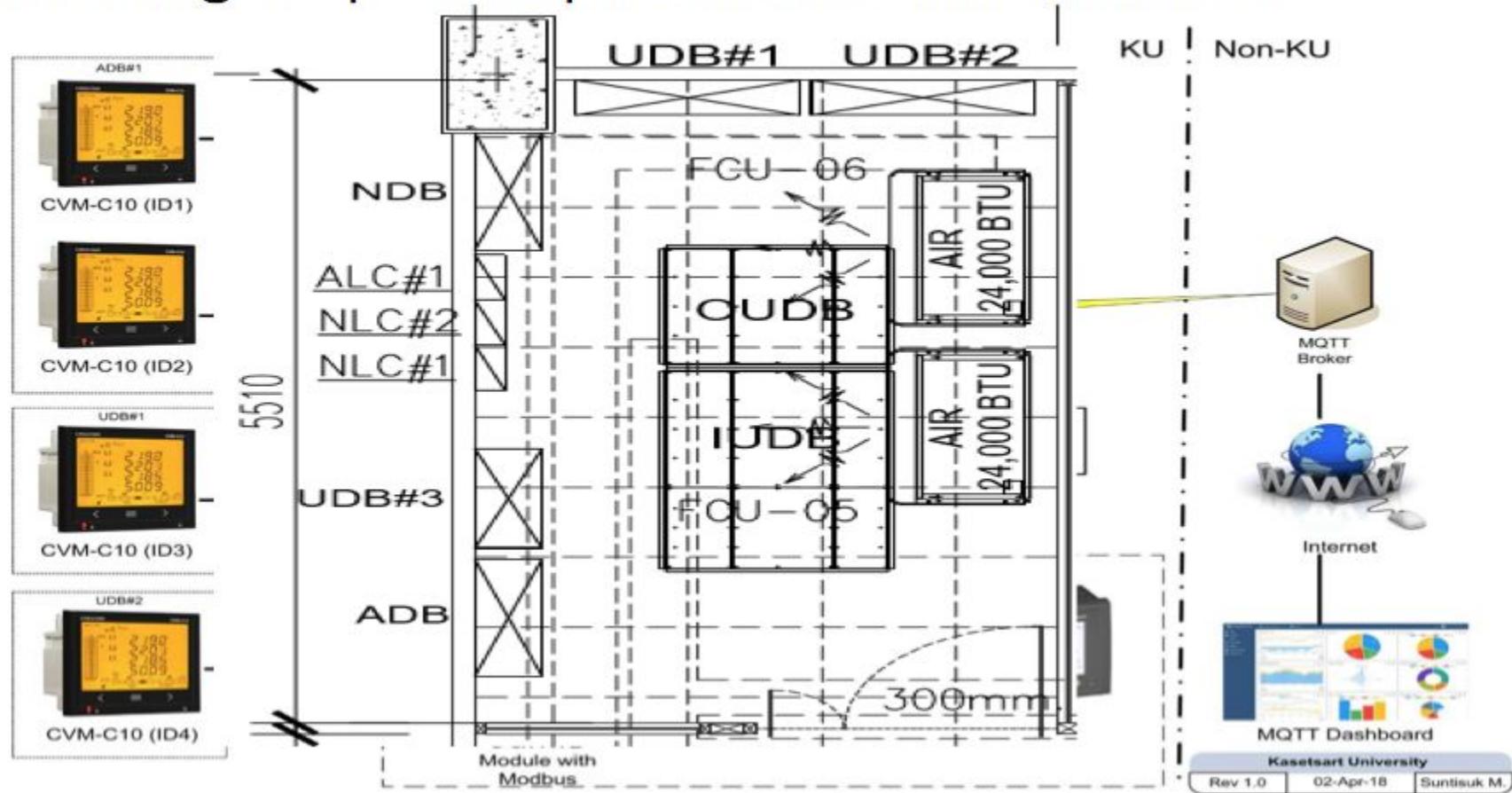
Kasetsart University



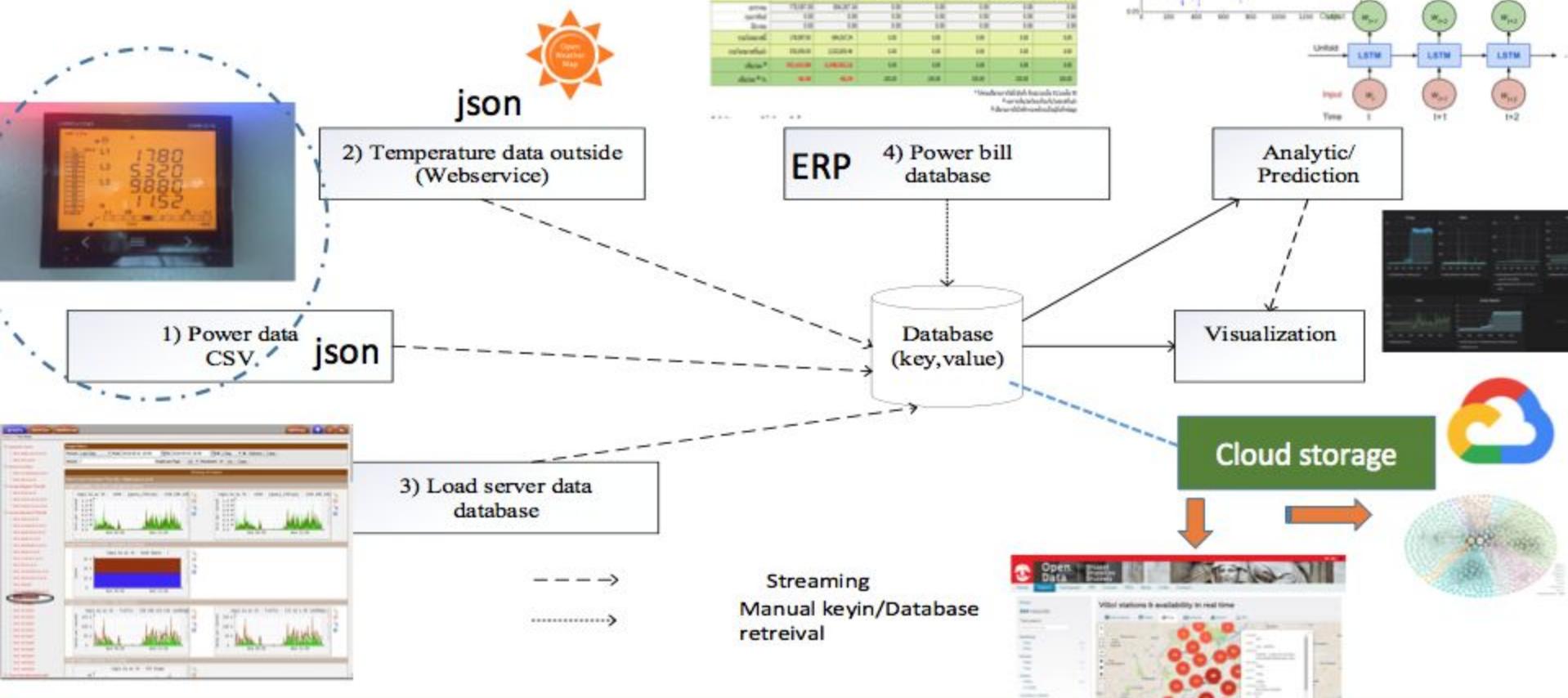


แบบแปลนค่าไฟฟ้าและครุภัณฑ์ สำหรับปรับปรุง
รวมทั้งหมด 1 : 100 ล้านเมตร

Measuring scope and power meter installation.



Relevant data integration



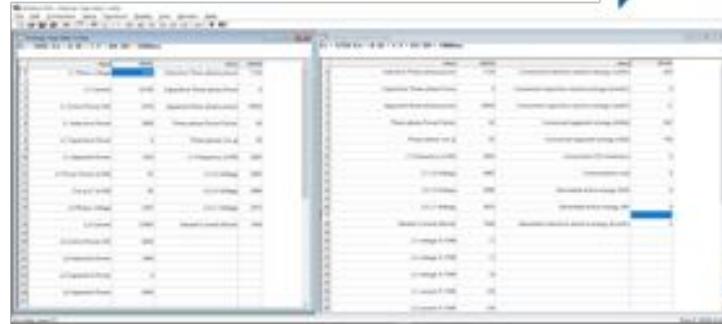
Meter installation

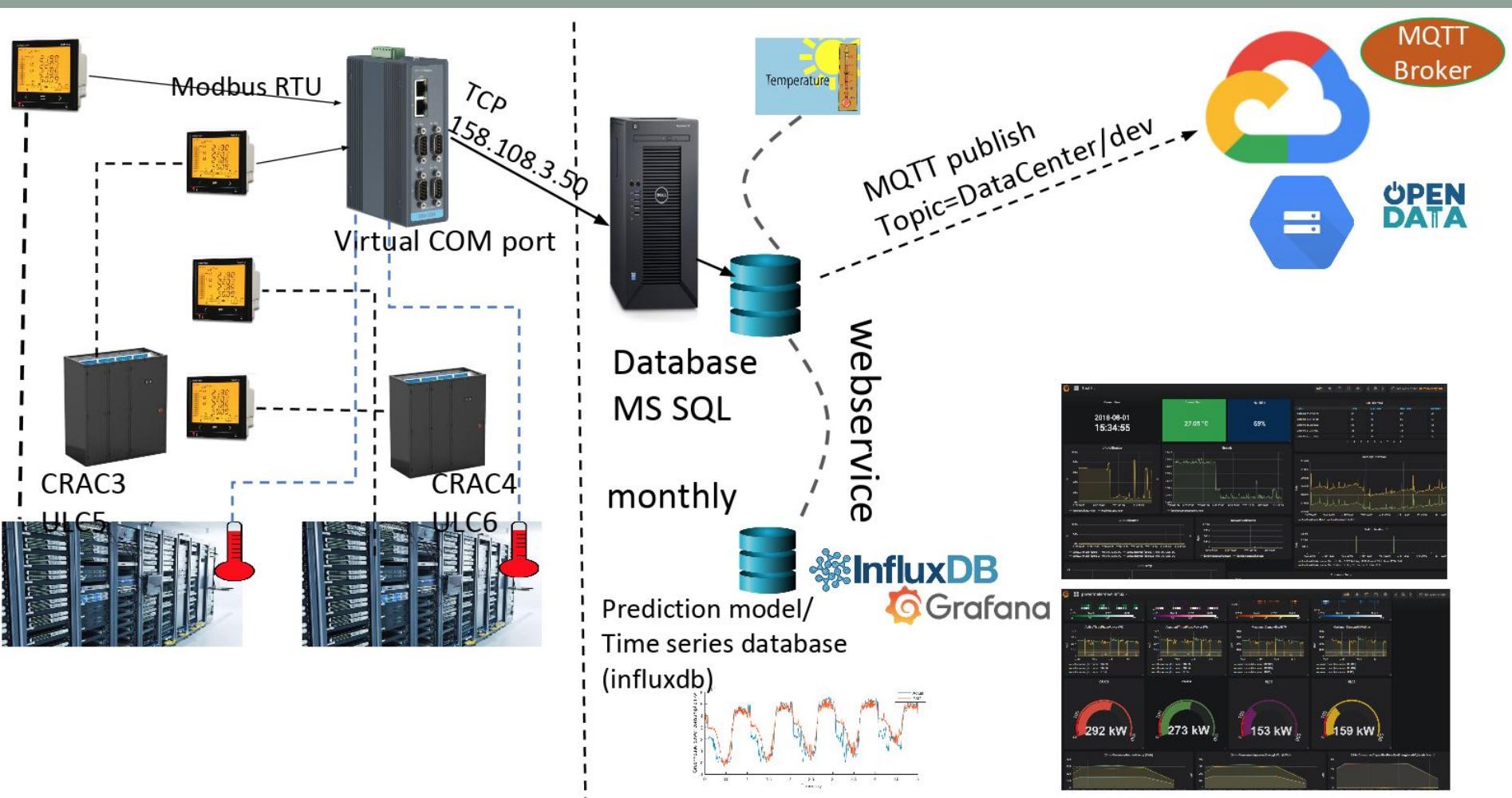


 CIRCUTOR



 NATIONAL
INSTRUMENTS





Visualization

UDB6 power today



UDB1_ULC5 today



CRAC3 power today



CRAC4 Power today



power meter sum -



Zoom Out



May 16, 2018 10:01:22 to a few seconds ago



Consumed active energy



value.Consumed_active_energy_kw (dev_name: ADB1_CRAC3)
value.Consumed_active_energy_kw (dev_name: ADB1_CRAC4)
value.Consumed_active_energy_kw (dev_name: UDB1_ULC5)
value.Consumed_active_energy_kw (dev_name: UDB2_ULC6)

Active three phase power



value.mean (dev_name: ADB1_CRAC3) value.mean (dev_name: ADB1_CRAC4)
value.mean (dev_name: UDB1_ULC5) value.mean (dev_name: UDB2_ULC6)



(ConsumedActiveEnergykW_kWh+ConsumedActiveEnergyW_Wh)



ActiveThreePhasePower_W



ConsumedCapacitiveReactiveEnergyvarhC_varh



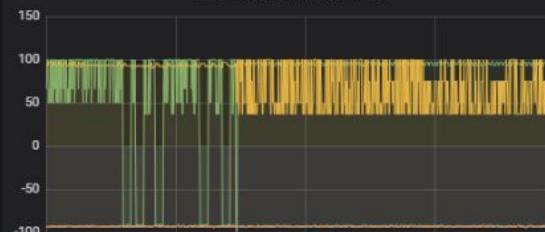
MaximumDemandkWIII_W



MaximumDemandkVAII



ThreePhasePowerFactorx100





Test0



Zoom Out

Today so far

Refresh every 30s



Current time

2018-05-25
10:11:18

Current Temp

26.00 °C

Humidity

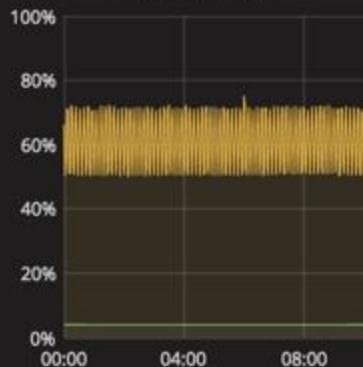
79%

Past Temp Value

Time	Temp	Min Temp	Max Temp	humidity
2018-05-25 09:51:42	31	30	31	79
2018-05-25 09:27:13	29	29	29	70
2018-05-25 08:51:42	27	27	28	83
2018-05-25 08:27:13	27	27	27	88
2018-05-25 07:51:42	27	27	27	88

1 2 3 4 5 6 7 8 9

CPU Utilization



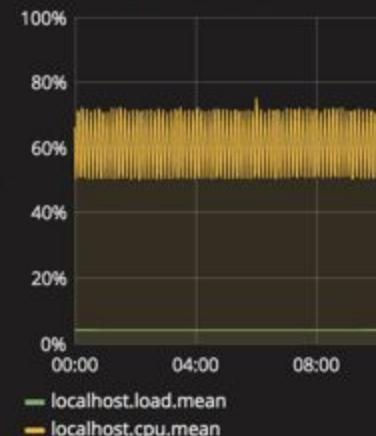
Threads



Memory Utilization



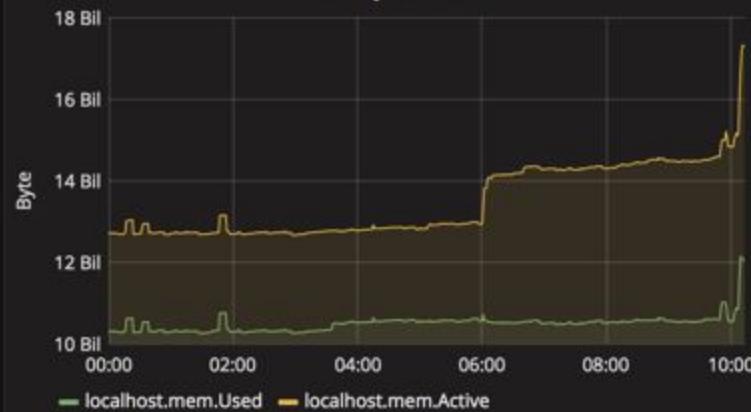
CPU Utilization



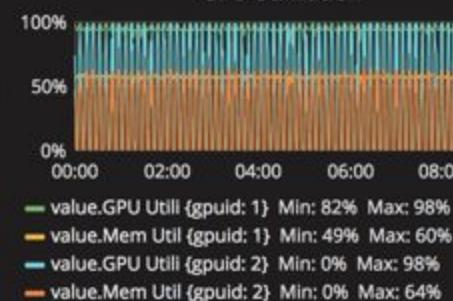
Threads



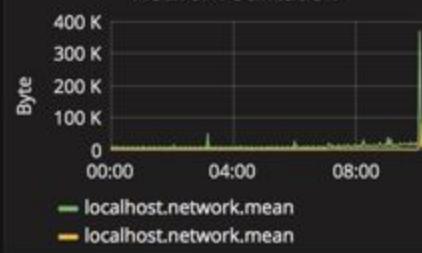
Memory Utilization



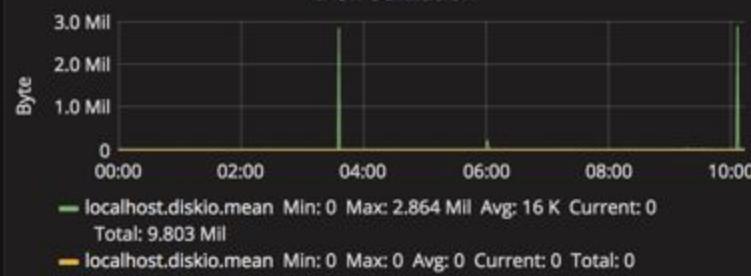
GPU Utilization



Network Utilization



Disk Utilization



Let's try out the data

Data Science Process

Data cleansing

Finding relationship

Modeling

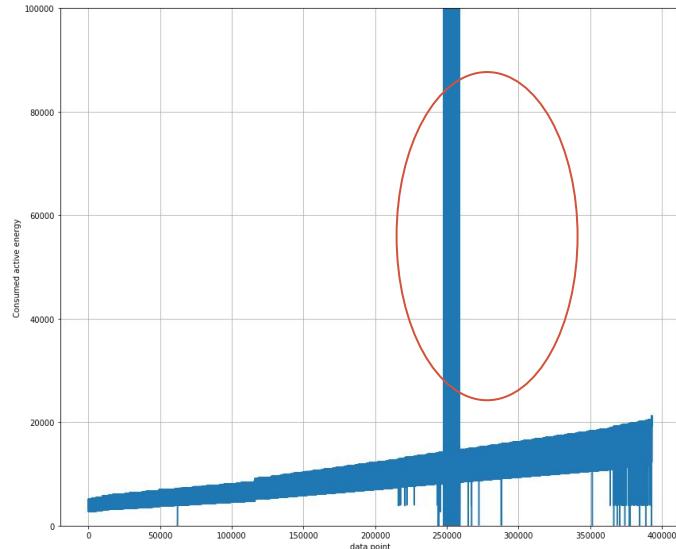
Checking accuracy

Visualization

Cleansing

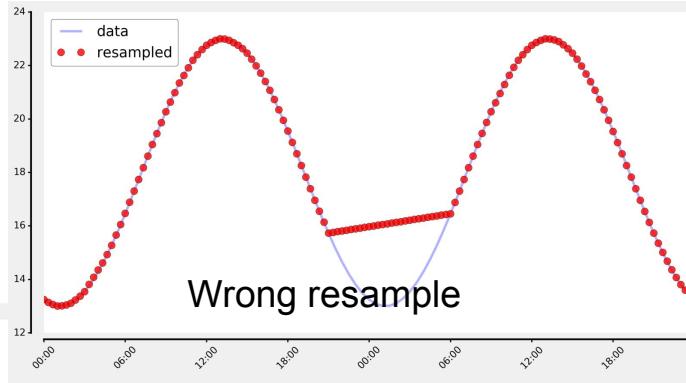
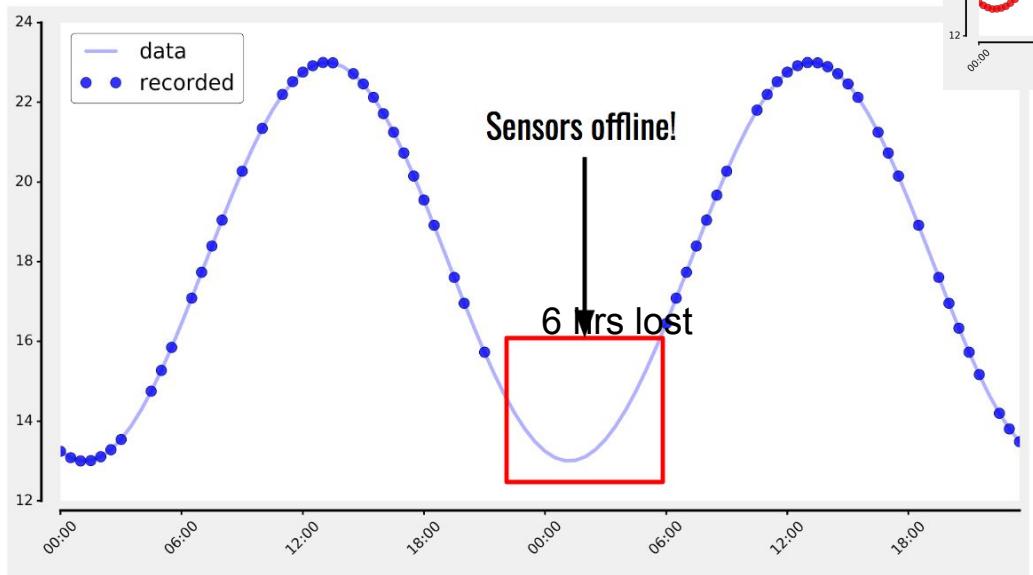
Remove noises

- Constraints (Max,Min)
- monotonically increasing



Missing data with resampling/interpolation

- Upsampling : Min -> Second
- Downsampling: Min -> Year



original

Consumed_active_energy_kW Consumed_apparent_energy_kVAh \

Timestamp

2018-05-16 14:41:16	5180.0	5413.0
2018-05-16 14:41:46	5180.0	5413.0
2018-05-16 14:42:16	5180.0	5413.0
2018-05-16 14:42:46	5180.0	5413.0

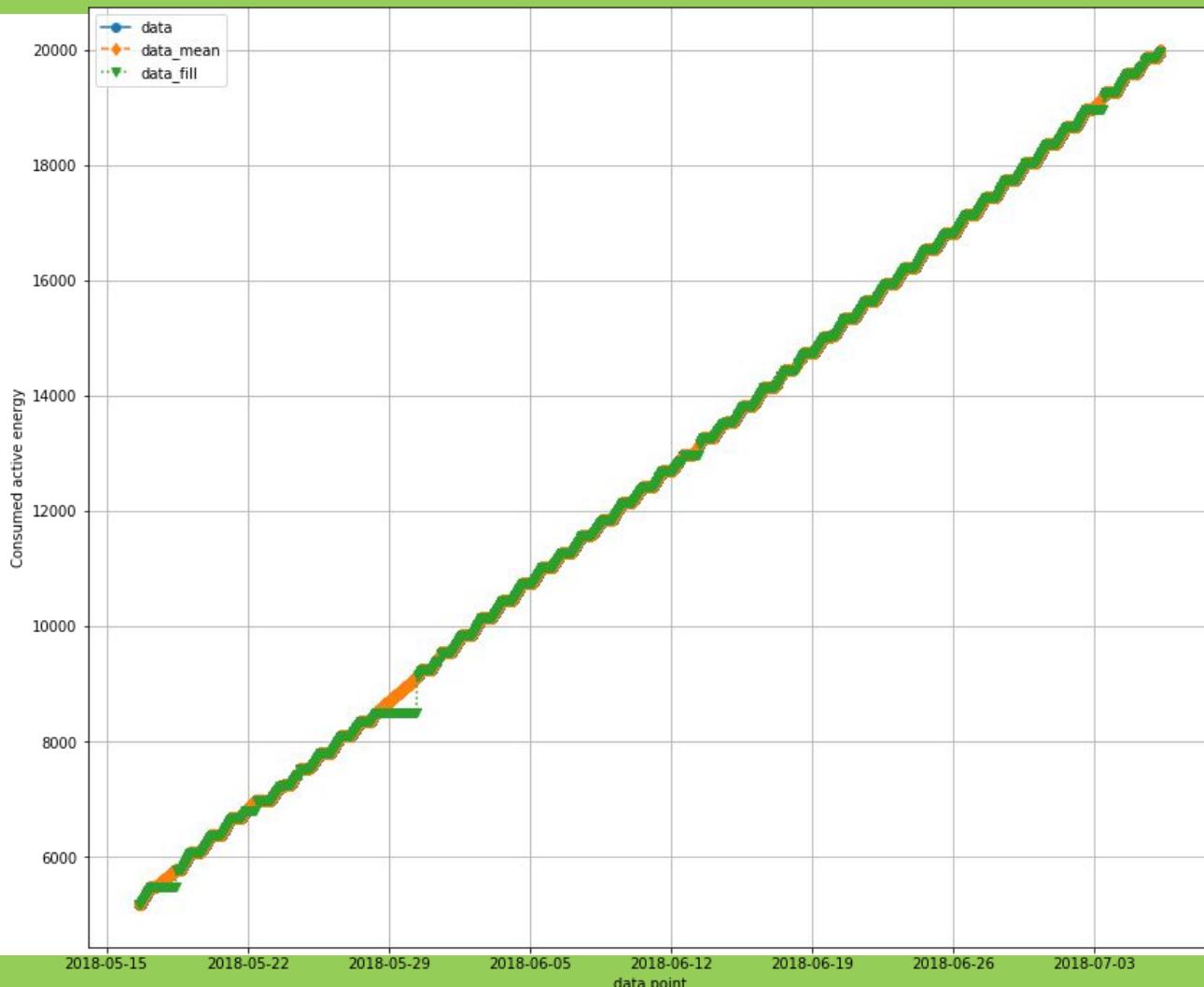
Every 2 m

2018-07-06 05:48:00	5254.000000
2018-07-06 05:50:00	5254.000000
2018-07-06 05:52:00	5254.000000
2018-07-06 05:54:00	5255.000000

[https://pandas.pydata.org/pandas-d
ocs/stable/generated/pandas.DataFrame.
resample.html](https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.resample.html)

Every 50min

2018-05-16 17:30:00	5205.0	5439.0
2018-05-16 18:20:00	5227.0	5461.0
2018-05-16 19:10:00	5248.0	5483.0
2018-05-16 20:00:00	5269.0	5504.0
2018-05-16 20:50:00	5290.0	5526.0
2018-05-16 21:40:00	5312.0	5549.0



Finding relationships with existing attributes

Some analysis...

1. Multivariable linear regression

Finding correlation

It is a mutual relationship between quantities.

- Sales increases when increase marketing budget?
 - Energy usages increase when temperature is high?
 - Customer purchases more if there are packages promotions?
- etc.

Intro to correlation

It is a mutual relationship between quantities.

- Sales increases when increase marketing budget?
 - Energy usages increase when temperature is high?
 - Customer purchases more if there are packages promotions?
- etc.

Why we need correlation?

Help predict one quantity from another.

Indicate causal relationship.

Basic foundation to other modeling techniques.

Type of correlations

Covariances:

Covariance is a statistical measure of association between two variables X and Y .

Indicate the increase in the same direction.

$$Cov(X, Y) = E[(X - E[X])(Y - E[Y])]$$

if both variables tend to move in the same direction, we expect the "average" rectangle connecting each point (X_i, Y_i) to the means (\bar{X}, \bar{Y}) to have a large and positive diagonal vector, **corresponding to a larger positive product in the equation above.**

Pearson Correlation Coefficient

Pearson correlation measures the linear association between continuous variables.

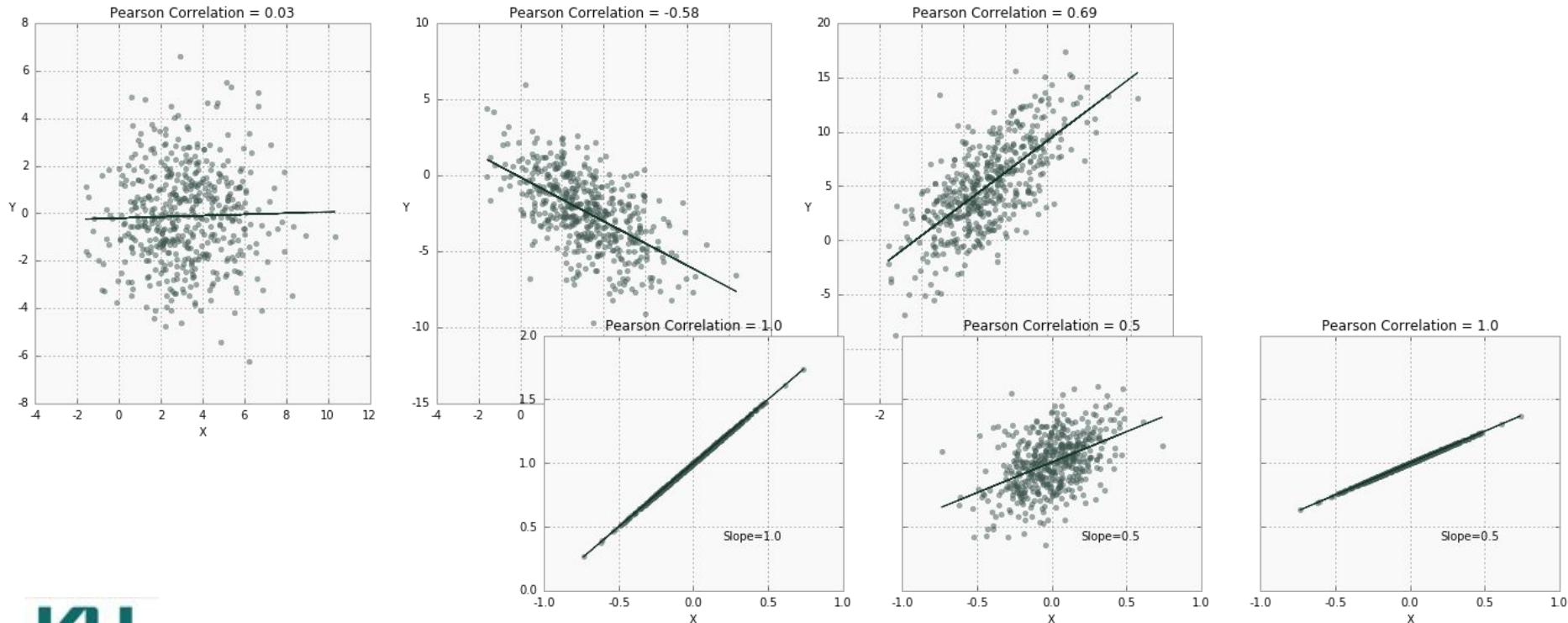
This coefficient quantifies the degree to which a relationship between two variables can be described by a line.

Developed by Karl Pearson over 120 years ago.

$$\rho_{X,Y} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$$

-Note that if X and Y are independent, then ρ is close to 0, but not vice versa!

-correlation does not equal slope



Spearman's Correlation Coefficient

- Spearman's correlation is not restricted to linear relationships.
- it measures **monotonic association** (only strictly increasing or decreasing, but not mixed) between two variables and relies on **the rank order of values**.

Spearman's coefficient looks at the relative order of values for each variable.

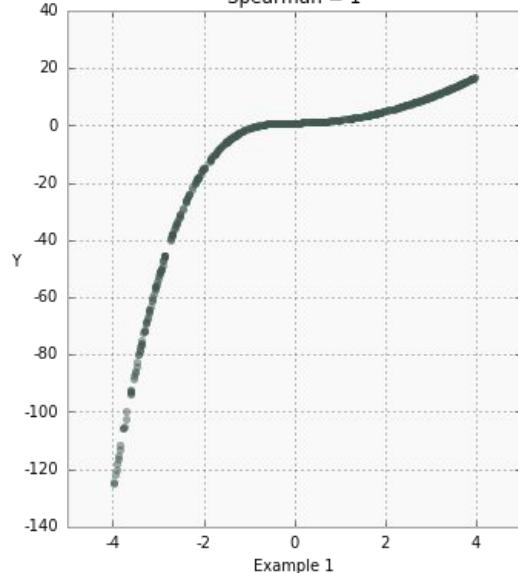
$$\rho_{rank_X, rank_Y} = \frac{cov(rank_X, rank_Y)}{\sigma_{rank_X} \sigma_{rank_Y}}$$

If all ranks are unique (i.e. there are no ties in ranks),

The simplified version: $\rho_s = 1 - \frac{6 \sum d_i^2}{N(N^2 - 1)}$

where $d_i = rank(X_i) - rank(Y_i)$,is the difference between the two ranks of each observation and N is the number of observations.

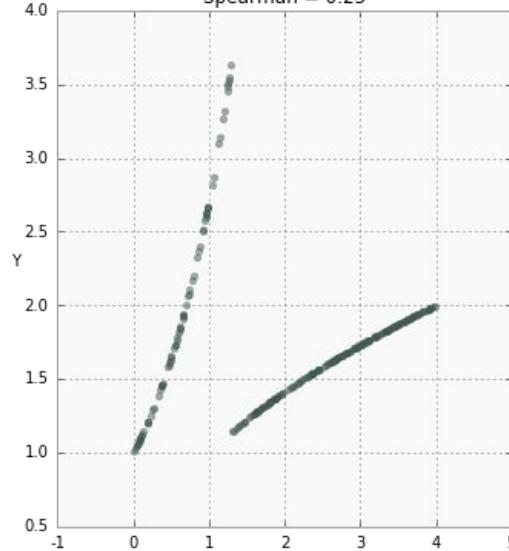
Pearson = 0.80,
Spearman = 1



Example 1

a clear monotonic (always increasing) and non-linear relationship.

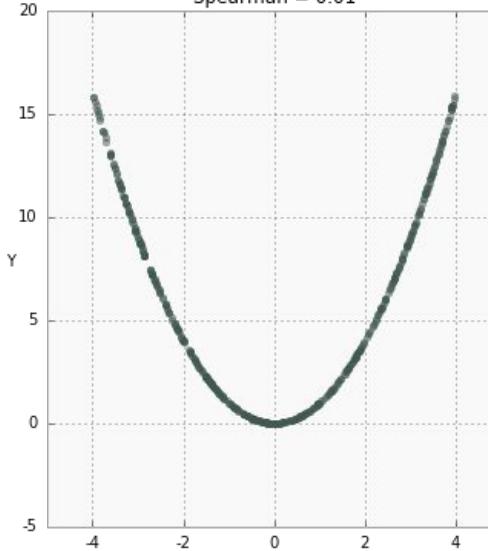
Pearson = 0.01,
Spearman = 0.25



Example 2

clear groups in X and a strong, although non-monotonic, association for both groups with Y . Pearson correlation is almost 0 since the data is very non-linear. Spearman rank correlation shows a weak association since the data is non-monotonic.

Pearson = 0.02,
Spearman = 0.01



Example 3

a nearly perfect quadratic relationship centered around 0. However, both correlation coefficients are almost 0 due to the non-monotonic, non-linear, and symmetric nature of the data.

Kendall's Tau Coefficient

Kendall's τ does not take into account the difference between ranks — only directional agreement. Therefore, this coefficient is more appropriate for discrete data.

$$\tau = \frac{(\text{number of concordant pairs}) - (\text{number of discordant pairs})}{N(N - 1) / 2}$$

Concordant pairs $(x_1, y_1), (x_2, y_2)$ are pairs of values **in which ranks coincide**: $x_1 < x_2$ and $y_1 < y_2$ or $x_1 > x_2$ and $y_1 > y_2$.

	X	Y
a	1	7
b	2	5
c	3	1
d	4	6
e	5	9

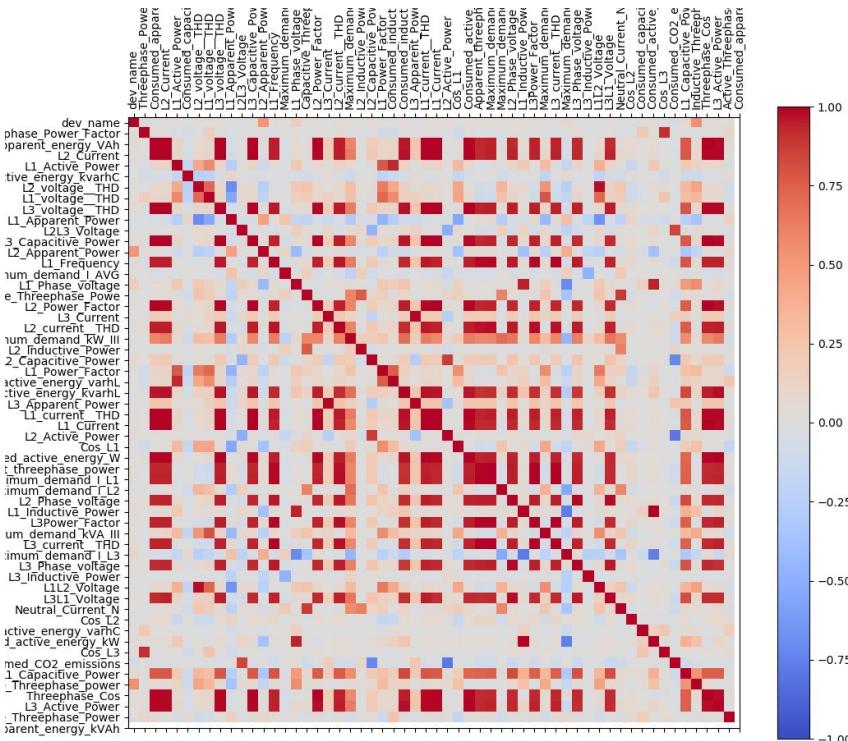
Here (4,6) and (5,9) in rows d and e is a concordant pair.
To calculate the numerator of τ , we compare all possible pairs in
the dataset and count number of concordant pairs; 6 in this case:

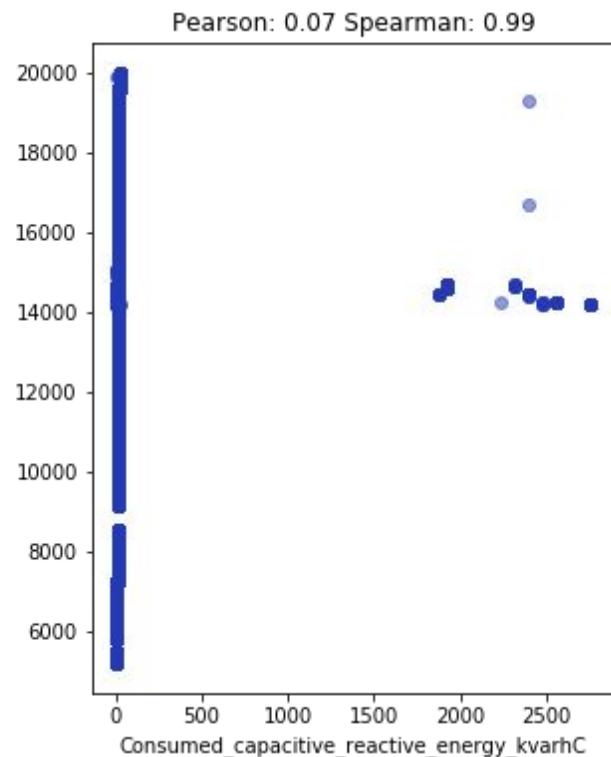
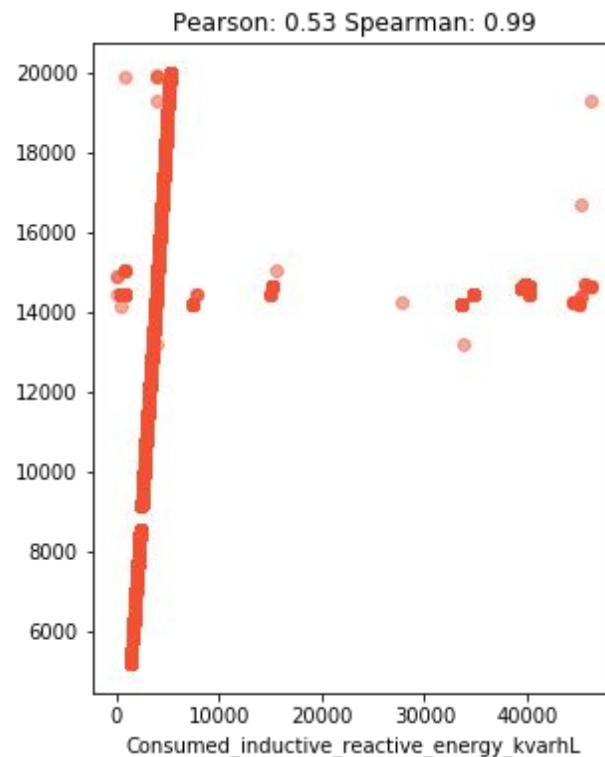
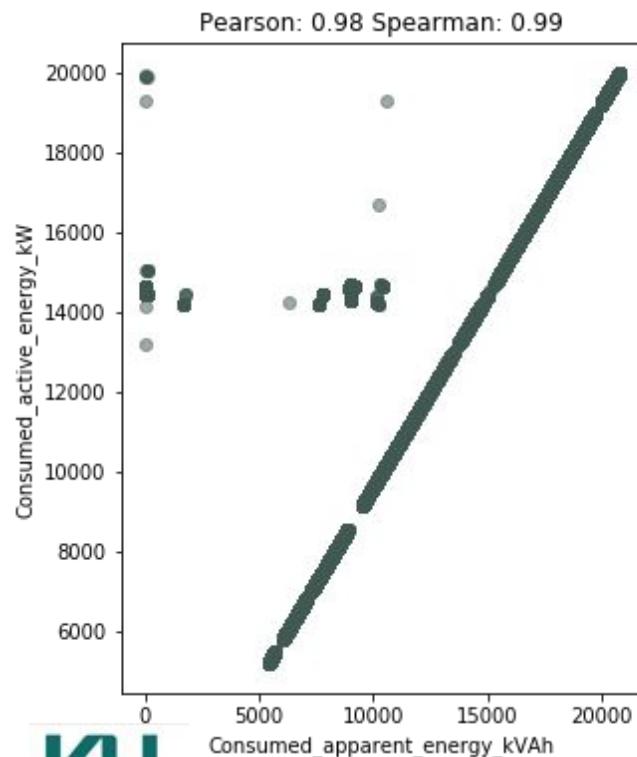
- (1,7) and (5,9)
- (2,5) and (4,6)
- (2,5) and (5,9)
- (3,1) and (4,6)
- (3,1) and (5,9)
- (4,6) and (5,9)

Discordant pair

- (1,7) and (2,5)
- (1,7) and (3,1)
- (1,7) and (4,6)
- (2,5) and (3,1)

Visualization of relation

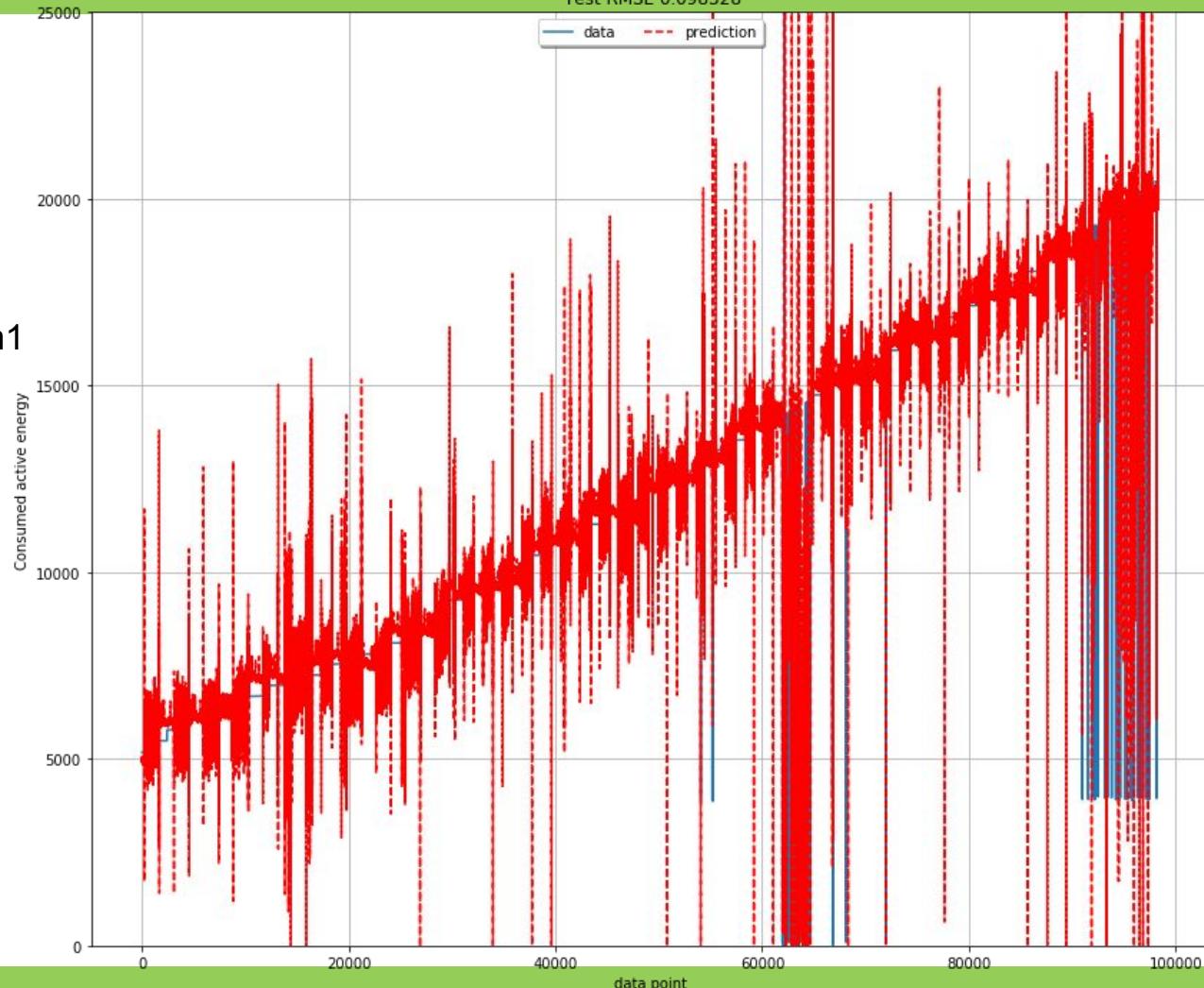




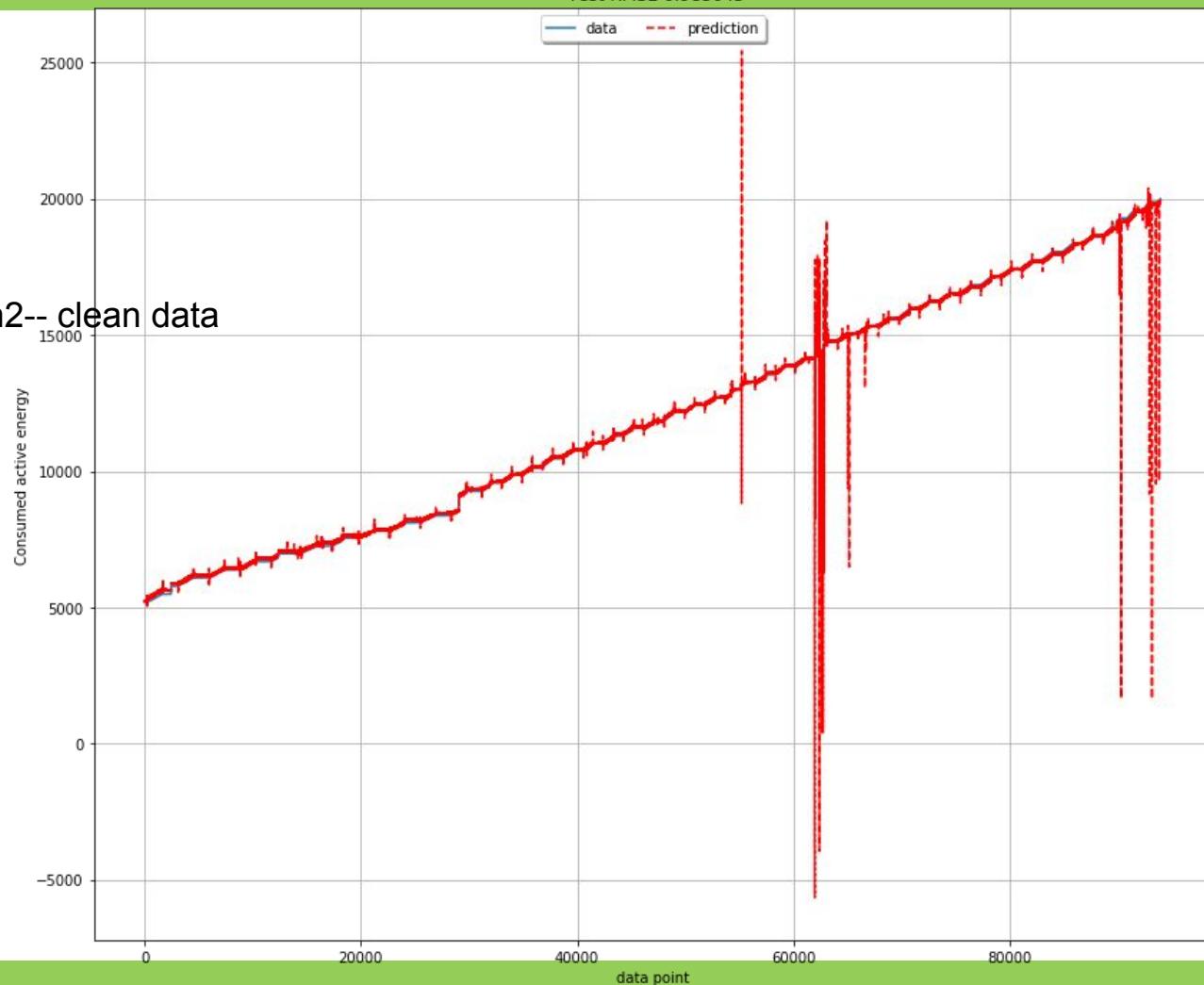
Top 20 features from correlation

Consumed_apparent_energy_kVAh	0.981747
Consumed_inductive_reactive_energy_kvarhL	0.528266
Consumed_capacitive_reactive_energy_kvarhC	0.068560
L2_voltage__THD	0.067717
Cos_L1	0.052542
L2_current__THD	0.050298
L1_current__THD	0.047905
Consumed_CO2_emissions	0.039708
Consumed_capacitive_reactive_energy_varhC	0.039587
L1_Power_Factor	0.035377
Threephase_Cos	0.034699
L1_Capacitive_Power	0.031062
L3_Phase_voltage	0.027613
Capacitive_Threephase_Powe	0.026782
Cos_L2	0.026488
Cos_L3	0.024519
L3Power_Factor	0.024410
L3_Apparent_Power	0.023013
Maximum_demand_I_L3	0.021808
L2_Power_Factor	0.021571

Test RMSE 0.098328



Test RMSE 0.985045

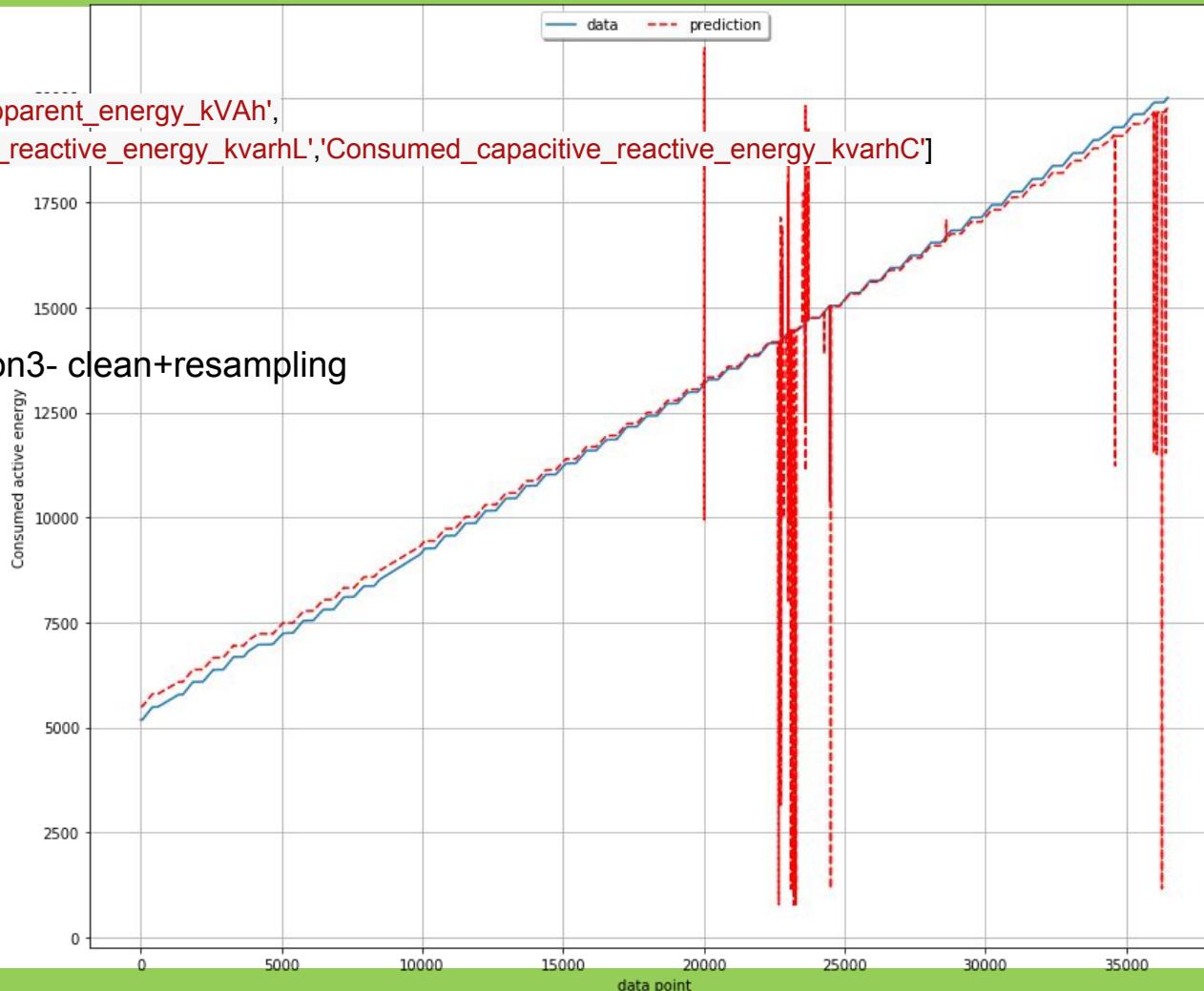


Test RMSE 0.956710

data prediction

```
cols = ['Consumed_apparent_energy_kVAh',  
'Consumed_inductive_reactive_energy_kvarhL','Consumed_capacitive_reactive_energy_kvarhC']
```

Linear regression3- clean+resampling



Resampling with
multicolumn



2. Autoregression

Autoregression

Time series data

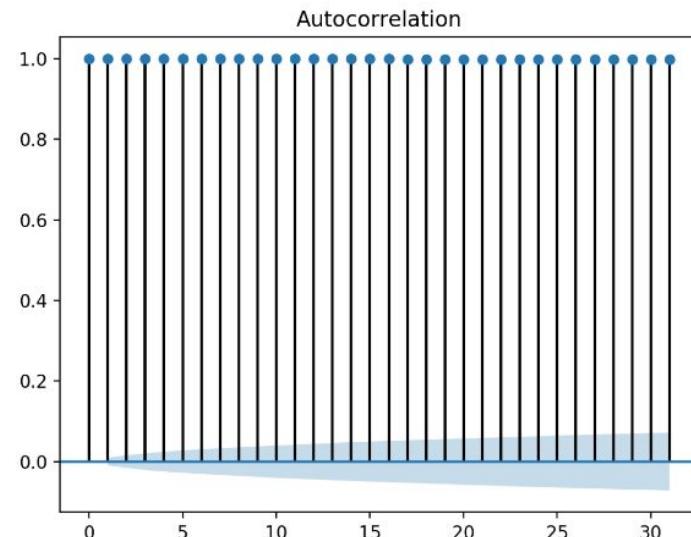
Predict $x(t+1)$ from $x(t-2)$, $x(t-1)$, $x(t)$

--Finding **autocorrelation**

Autocorrelation

	$X(t-1)$	$X(t+1)$
$X(t-1)$	1.0	1.0
$X(t+1)$	1.0	1.0

good!



Lag=31

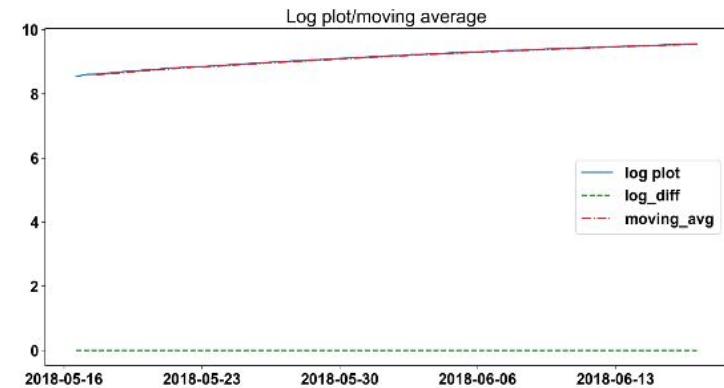
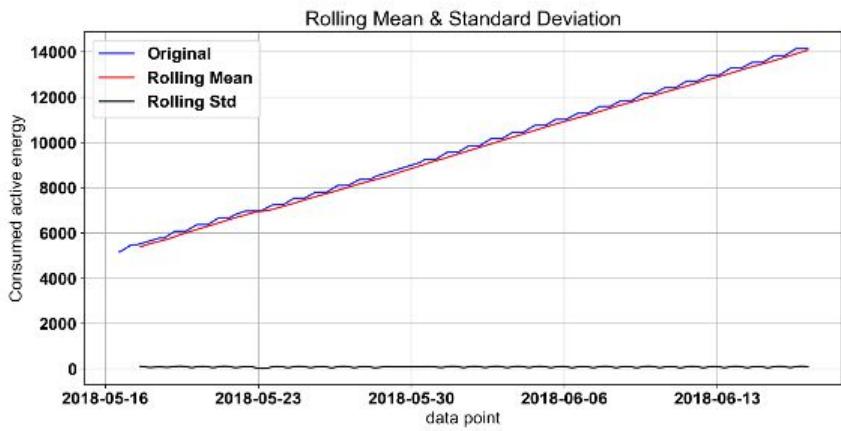
Stationary analysis

- Testing whether the data is stationary. It has the constant average/moving average
- Use Dickey Fuller test

Test statistics -0.061172 should be less than Critical Value (1%)

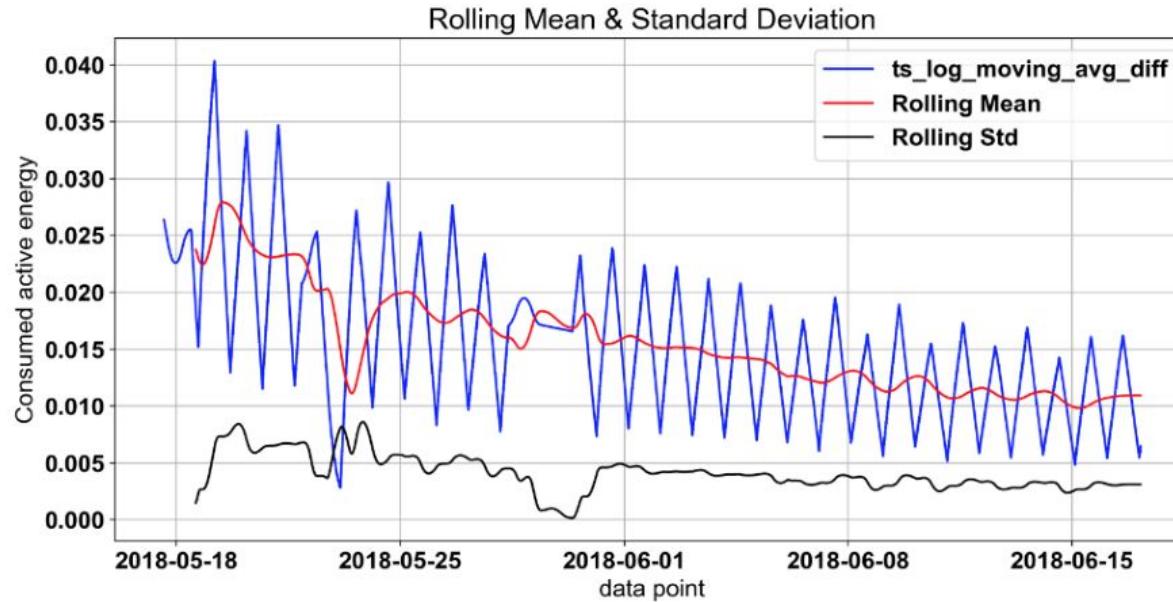
Results of Dickey-Fuller Test:

Test Statistic	-0.061172
p-value	0.953171
#Lags Used	55.000000
Number of Observations Used	45332.000000
Critical Value (5%)	-2.861604
Critical Value (1%)	-3.430494
Critical Value (10%)	-2.566804



Solving non-stationary

-Take difference (first order diff : with prev value)

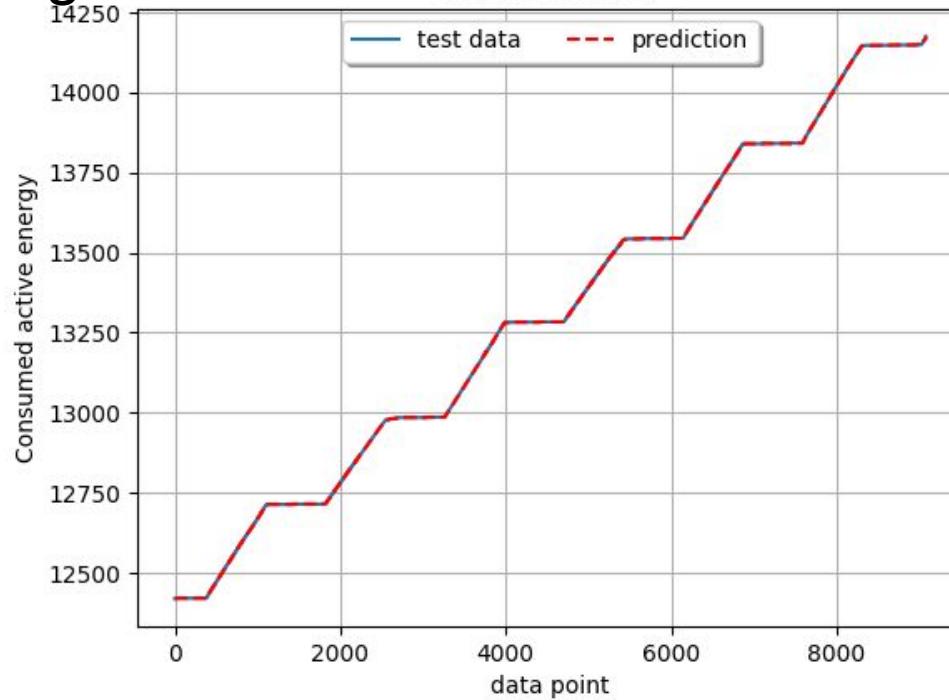


Test Statistic **-8.922295e+00**
p-value 1.033276e-14
#Lags Used 5.500000e+01
Number of Observations Used 4.389300e+04
Critical Value (5%) **-2.861606e+00**
Critical Value (1%) **-3.430499e+00**
Critical Value (10%) -2.566805e+00

Better!!

Prediction with autoregression

RMSE 0.393294



<https://cocalc.com/projects/>

Free Jupyter project: jupyter lab

https://github.com/cchantra/datascience_tutorial.git

KU Data Science Status

KU Data Science Bootcamp

KU Data Science Forum

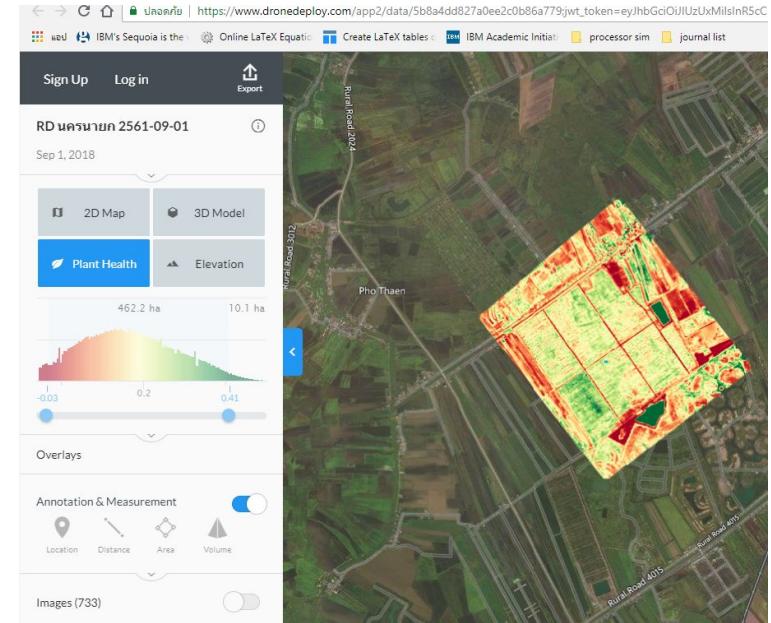
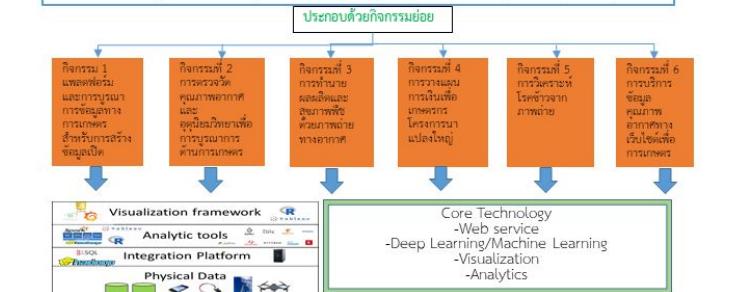
KU Data Science Education

KU Open data

<http://datascience.ku.ac.th>

วัตถุประสงค์

- สร้างระบบนำเสนอด้วยเทคโนโลยีทางเกษตรเพื่อการนำเสนอข้อมูลทางการเกษตรเพื่อส่งเสริมความเป็นไปได้ที่จะประดิษฐ์
- นำข้อมูลทางเกษตรที่ได้รับมาและรวมไว้ให้เข้ากันได้ที่ต้องการซึ่งเป็นข้อมูลเปิด (Open data) ทางภาคเกษตร
- สร้างแพลตฟอร์มที่สามารถใช้ข้อมูลเปิด (Open data) ทางภาคเกษตร
- สร้างบริการใช้ข้อมูลและวิเคราะห์ที่ช่วยเหลือให้เกิดขึ้น Machine learning และ Deep learning ได้





KU Open Data

ABOUT NEWS DATA COLLECTION DATA SUBMISSION

KU Open Data Collection

Browse by categories

Financial data

[Energy and Environment](#)

Registration

[Food & Agriculture](#)

Library

<https://www.facebook.com/KU-Data-science-boot-camp-2018-356005521545466/?ref=bookmarks>



KU Data Science Bootcamp 2018

20 May -26 May 2016



KU Data science
boot camp 2018

Create Page @Username

Home



Posts

Reviews

Videos

Photos

KASETSART
UNIVERSITY

Like Follow Share ...

Send Message

Create Post Live Event Offer Job

AUDS
Test



Write a post...



No Rating Yet

