



INTRODUCTION TO DATA ANALYTICS

Introduction to the Data Analytics Course

Asst. Prof. Dr. Rathachai Chawuthai

Department of Computer Engineering
Faculty of Engineering
King Mongkut's Institute of Technology Ladkrabang

Agenda

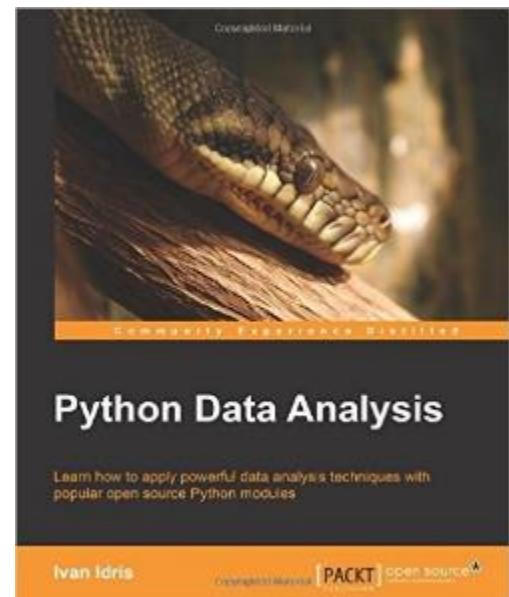
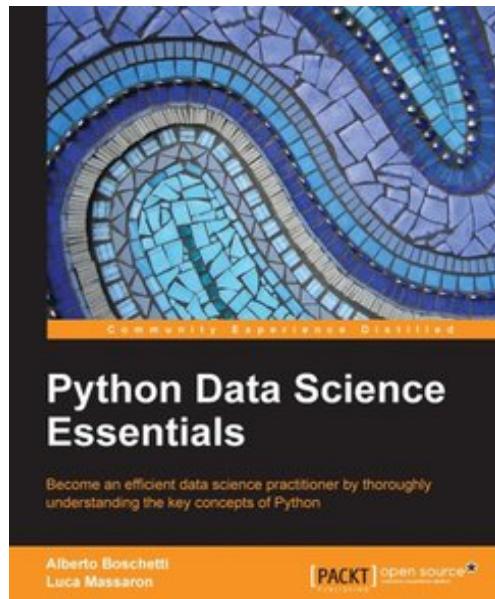
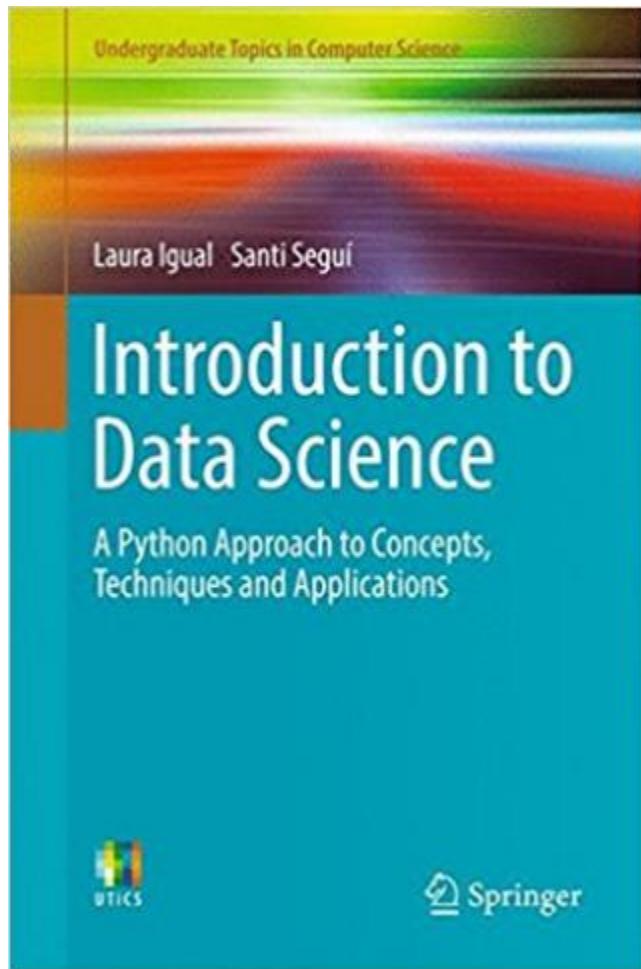
- 
- About this Course
 - Data Analytics

About this Course

Course Outline

- Introduction
- Python for DA
- Data Exploration
- Data Processing
- Regression Analysis
- Classification Analysis
- Cluster Analysis
- Recommender System
- Visualization

Books



[http://www.springer.com/
gp/book/9783319500164](http://www.springer.com/gp/book/9783319500164)

E-Learning

Python for Data Science and Machine Learning Bootcamp

Learn how to use NumPy, Pandas, Seaborn , Matplotlib , Plotly , Scikit-Learn , Machine Learning, Tensorflow , and more!

★★★★★ 4.5 (60,222 ratings) 285,702 students enrolled

Created by Jose Portilla Last updated 9/2019

English English, Indonesian [Auto-generated], 7 more



Preview this course

i You purchased this course on Oct. 7, 2017

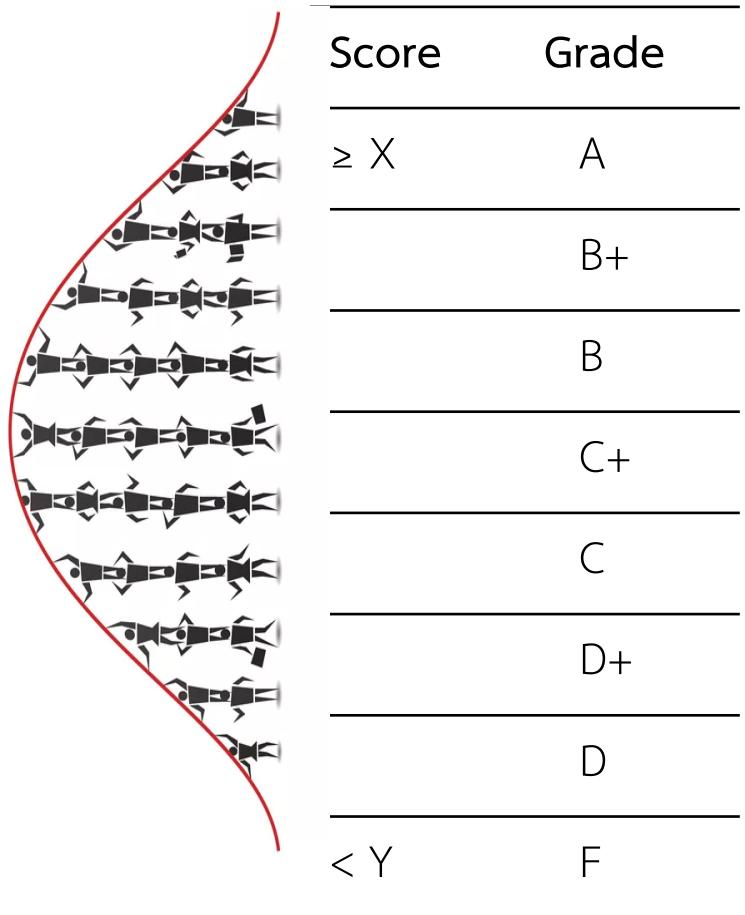
[Go to course](#)

<https://www.udemy.com/course/>

[python-for-data-science-and-machine-learning-bootcamp/](https://www.udemy.com/course/python-for-data-science-and-machine-learning-bootcamp/)

Grading (Plan)

- Midterm 30%
- Final 30%
- Assignments 10%
- Project 30%



Tools

Anaconda

- Language: Python
- Package: NumPy, SciPy, Pandas, Scikit-learn, , etc.
- IDE: Spyder, Jupyter
- Download:

<https://www.anaconda.com/download/>

Data Analytics

Questions ?

ទូរសព្ទរបាយការណ៍របស់អ្នក ត្រូវបានបញ្ជាក់ដោយចំណាំ

How much is the price of 1 acre of land in this area?

How to commute from here to Icon Siam at 8:00?

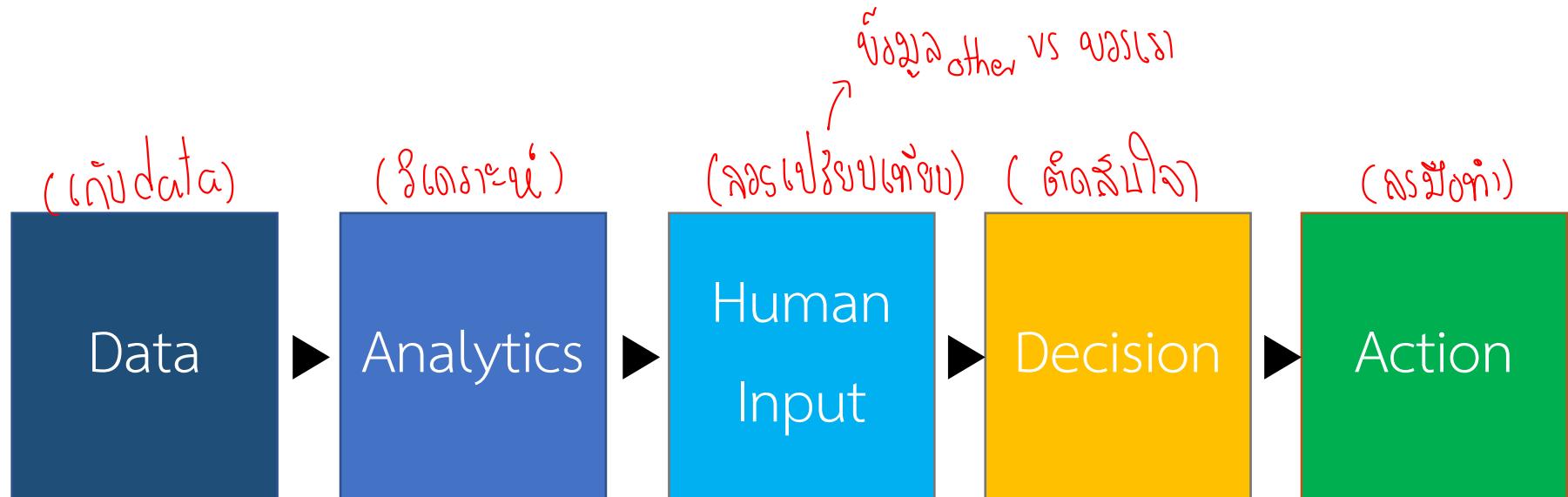
Will this be a flu or not?

If the customer has already bought this item,
which products should be recommended next?

Where should the gas station be opened?

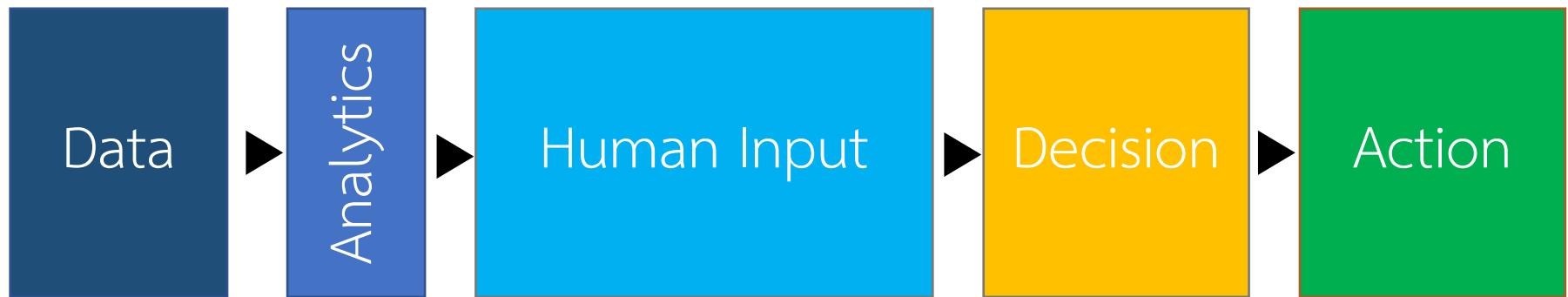
Data Analytics

ກວດສອບຕາມຂ່າຍ data ana

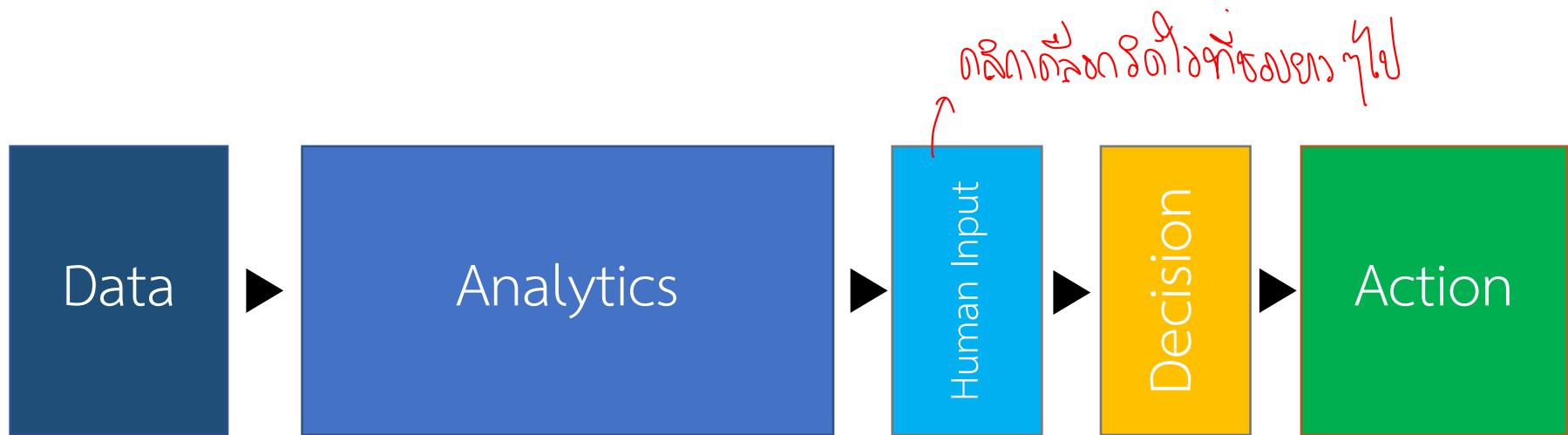


ດີ analytics ອື່ນ → ດົນເຕັ້ນທີ່ໄດ້ຢູ່ຂະໜາດ
ຕໍ່ການສຸດທະນາຖາວອນ ມານຸ່ງທີ່ໄວເຕັ້ນຮັດແກ່ໂຄດເລຍ

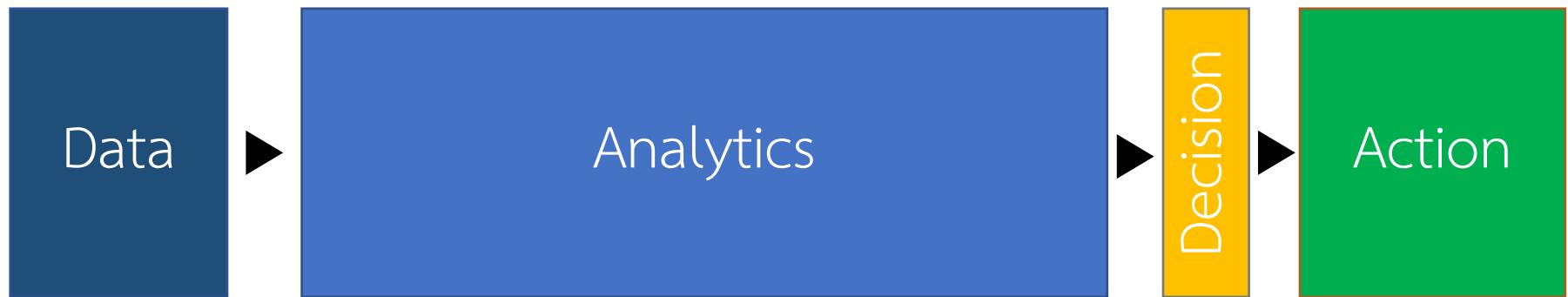
Data Analytics



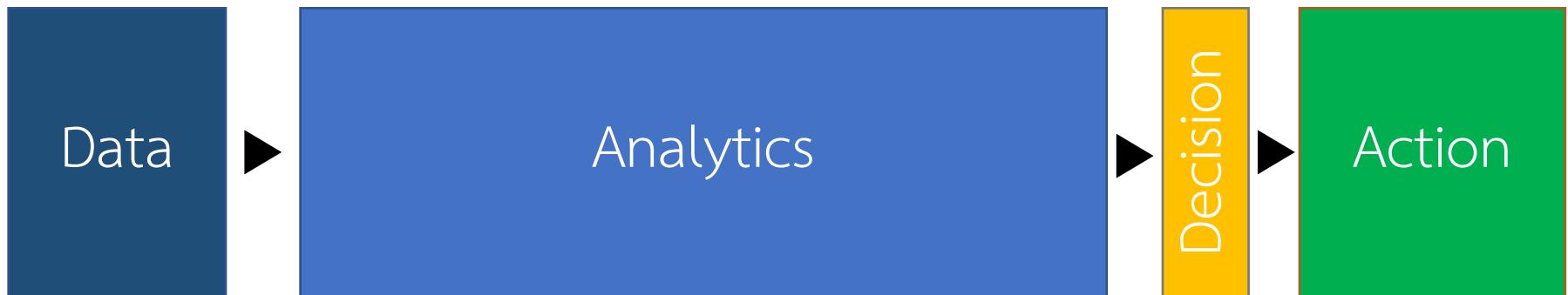
Data Analytics



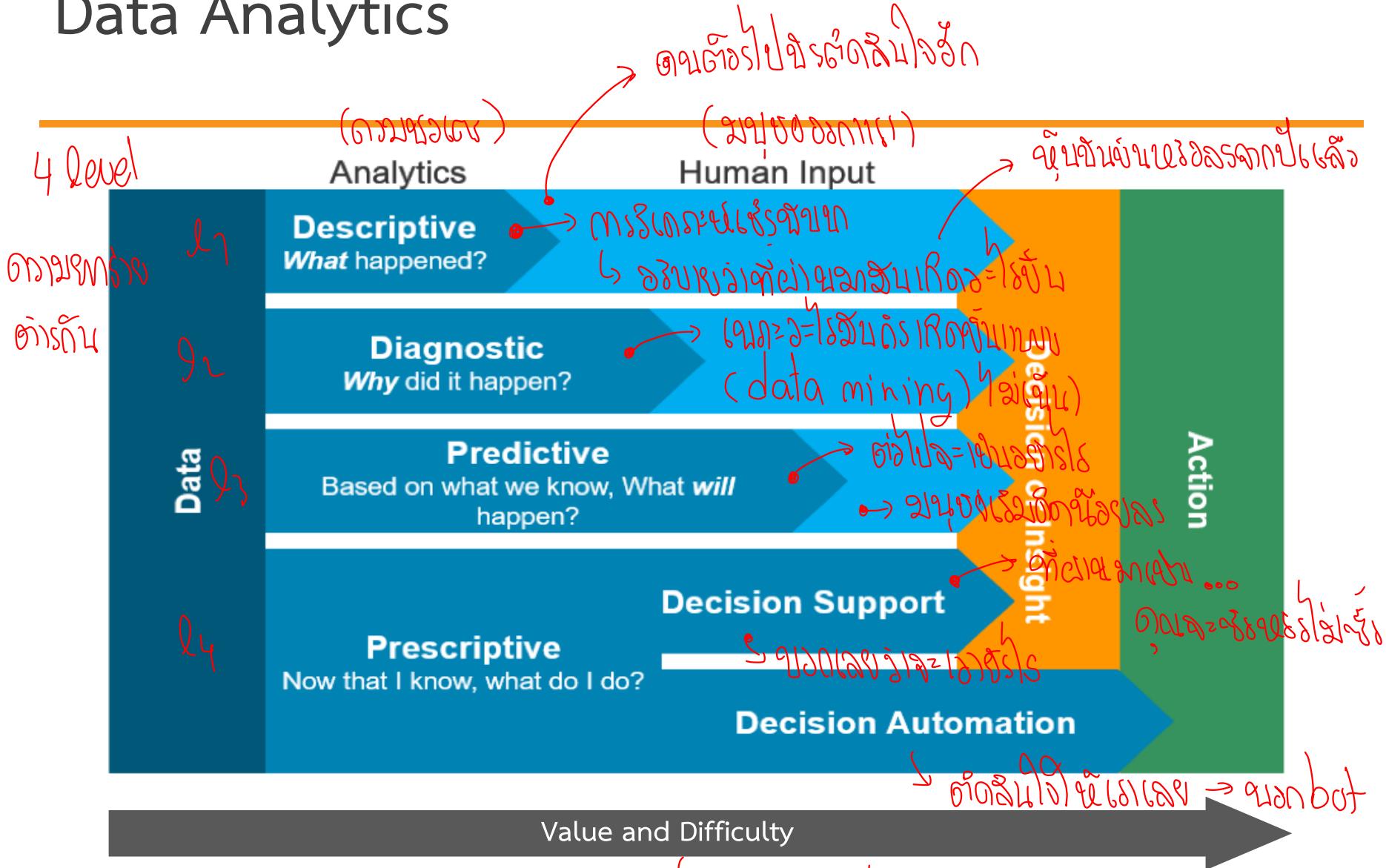
Data Analytics



Data Analytics



Data Analytics



- Ref:** • Four types of analytics capability (Gartner, 2014)
• (image) <https://www.healthcatalyst.com/closed-loop-analytics-method-healthcare-data-insights>

Descriptive Analytics → ຕິດການທີ່ມີນຳໄລ້ຕໍ່າງໆ → ສິນຕົວໄປທາຍຫຸ່ນໆ, ສັງເກດທີ່ມີການ

Descriptive analytics is the **interpretation of historical data** to better understand changes that have occurred in a business. Descriptive analytics describes the use of a range of historic data to draw comparisons.



designed by freepik

Descriptive Analytics

Questions:

- Which are the best-seller products?
- Which are the most or least revenue-generating products?
- Which are the most successful promotional campaigns?
- Who are the most paying customers?
- What are revenue trends for each Strategic Business Unit (SBU) of last N years, last N months?

Descriptive Analytics

Techniques:

ឧបករណ៍សារិកជាមុន

- Exploratory Data Analysis
- Measure of the Shape of the Distribution
- Measure of Data Summary
- Measure of Variability or Dispersion
 - Standard Deviation, Interquartile Range, Range
- Measure of Central Tendency
 - Mean, Median, Mode, Min, Max

Predictive Analytics → ඉංග්‍රීසුයාධාරිත → සිද්ධාන්ත

Predictive analytics is the practice of extracting information from existing data sets in order to **determine patterns and predict future outcomes and trends.**

Predictive analytics does not tell you what will happen in the future.



Predictive Analytics

business analytics
ការប្រើប្រាស់នូវការសម្រាប់
ត្រួតពិនិត្យនូវការងារដែលអាចរាយការណា

Questions:

របៀបរាយការណ៍ ត្រួតពិនិត្យនូវការ

- What is going to be likely revenue for each SBU in coming year?
- What is going to be likely attrition rate for the common year?
- Who all customers are likely to churn-out?
- Which promotional campaigns are likely to do well?
- Which products are likely to sell most in the next 6 months?

Predictive Analytics

Techniques:

- Decision Support System
 - Linear Regression
- Classification
 - Decision Tree
 - Logistic Regression
 - Support Vector Machine
 - Artificial Neural Network
 - etc.

Prescriptive Analytics

ក្រសួងអនុវត្តន៍របាយការ

Prescriptive analytics is a type of data analytics—the use of technology to help businesses make **better decisions** through the analysis of raw data. Specifically, prescriptive analytics factors information about possible situations or scenarios, available resources, past performance, and current performance, and suggests a course of action or strategy. It can be used to make decisions on any time horizon, from immediate to long term.



Prescriptive Analytics

→ វិនិច្ឆ័យ សមារករណីលើឈូលក់អុទ្ធតែងសំខាន់

តើពីរបៀវត្សនេះនឹងដឹងដីណា

Questions:

↓
តាមរយៈ A → រយៈ B [dif!]
រយៈ C → នូវរយៈ D

តុលាកំណែងដឹងដីណា

ការគិតថ្លែង project និងគោលរាល់ទេ

- What would be the best channel to sell this product?
- Which of the supplier suggested promotions of adopt?
- What new or replacement items to introduce, and when?
- How to modify the overall product assortment for each category?
- What's the next promotion that I can offer to this customer segment?
- What is the best route from the point A to B?

Prescriptive Analytics

Techniques:

- Decision Support System
- Recommender System
- Search Engine
- Route and Direction Recommendation
- Chatbot



??? Analytics ???

- What was the popular product last month? → descriptive
- What is the average revenue of this product? → descriptive
- What will be the revenue of the next quarter? → predictive
- Which products should be promoted next month? → prescriptive
- Which products should be stopped selling? → prescriptive
- Which place should we promote this product? → prescriptive
- Which products will we recommended to our customers? → prescriptive
- Will the customers cancel their orders? → prediction

When you have a question!

Finding Data

Analyzing

Finding Answer

OK?

Example: a small dataset

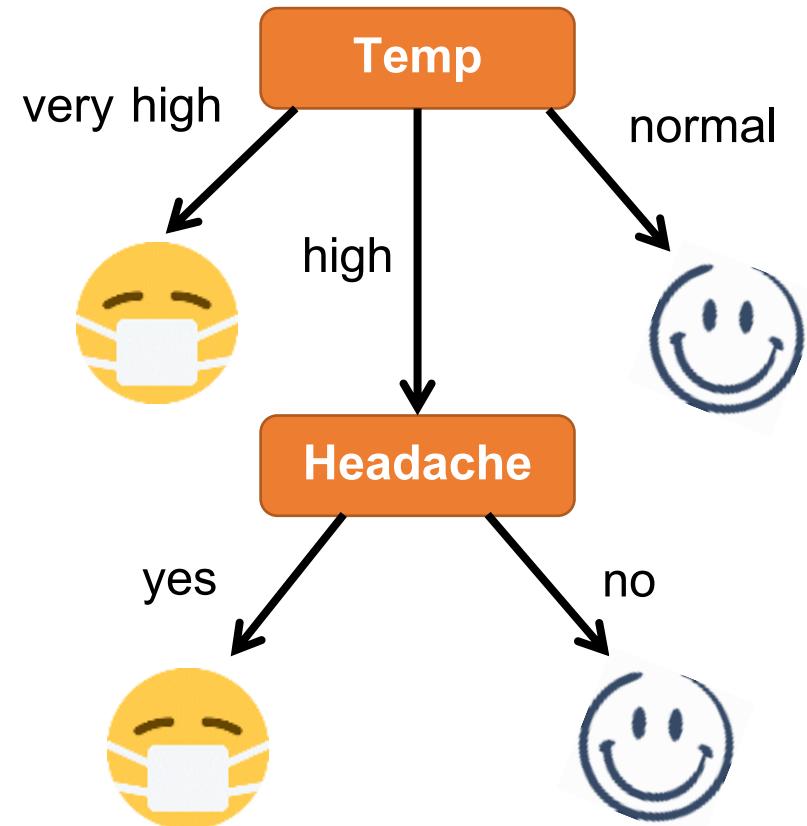
#	Body Temperature	Headache	Nausea	FLU?
1	high	yes	-	yes
2	very high	yes	yes	yes
3	normal	-	-	-
4	high	yes	yes	yes
5	high	-	yes	-
6	normal	yes	-	-
7	normal	-	yes	-

	Normal	YES	YES	?
--	--------	-----	-----	---

Example: a small dataset

#	Body Temp	Headache	Nausea	FLU?
1	high	yes	-	yes
2	very high	yes	yes	yes
3	normal	-	-	-
4	high	yes	yes	yes
5	high	-	yes	-
6	normal	yes	-	-
7	normal	-	yes	-

	Normal	YES	YES	?
--	--------	-----	-----	---



In this age

Much Data → High Accuracy

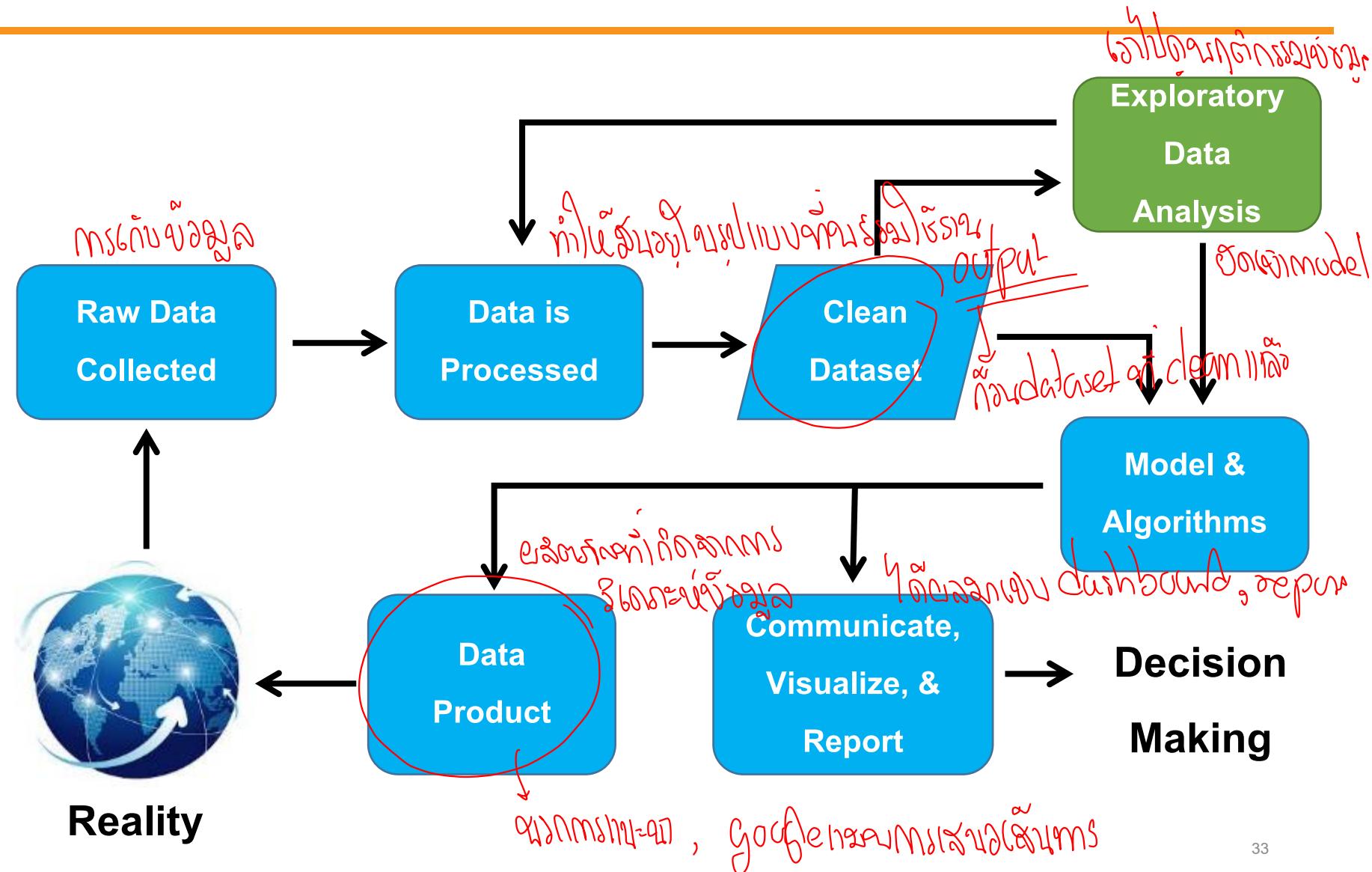
Much Data → Headache

(Then, let's use Computers)

Data Science

Data science, also known as data-driven science, is an interdisciplinary field about **scientific methods, processes and systems to extract knowledge or insights from data** in various forms, either structured or unstructured, similar to Knowledge Discovery in Databases.

Data Science Process



Weather Forecast

អេឡិចត្រូនកម្មណ៍ (data product)

Shibuya, Tokyo, Japan

Wednesday 6:00 AM

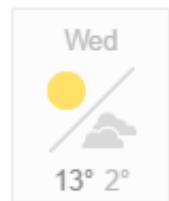
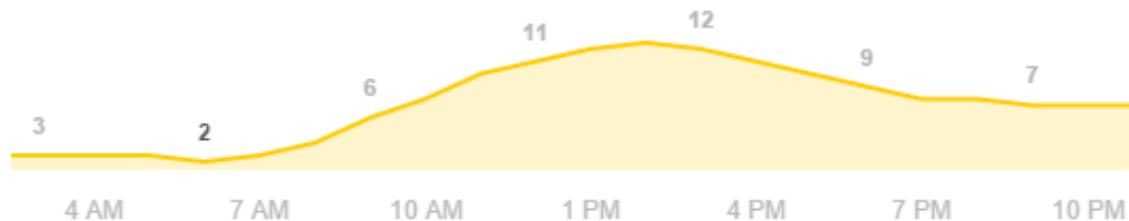
Clear



2
°C | °F

Precipitation: 0%
Humidity: 83%
Wind: 6 km/h

Temperature Precipitation Wind



amazon

Frequently Bought Together



Total price: To see our price, add these items to your cart. Why don't we show the price?

[Add both to Cart](#)
[Add both to List](#)

These items are shipped from and sold by different sellers. Show details

This item: Sony KDL40W650D 40-Inch 1080p Smart LED TV (2016 Model)

Cheetah Mounts APTMM2B TV Wall Mount for 20-75-inch TVs Bundle with 10-feet Braided HDMI Cable and a... \$24.99

Customers Who Bought This Item Also Bought

Page 1 of 7



Ultra High Speed HDMI Cable 1080p Cable for HDTV, Blu-Ray, PS3 (6 feet)
 340
\$12.20



VideoSecu ML531BE TV Wall Mount for most 22"-55" LED LCD Plasma Flat Screen Monitor up to 88 lb VESA 400x400mm
 17,844



Articulating Arm 32-50 inch TV LCD Monitor Wall Mount, Full Motion Tilt Swivel for 32" 36" 37" Flat Screen Monitor up to 88 lb VESA 400x400mm
 170



VideoSecu TV Wall Mount Tilt Low Profile Ultra Slim Television Mount Bracket for Most 26"- 47" LED LCD Flat Screen Monitors up to 88 lb VESA 400x400mm
 588



Sony XBR49X800D 49-Inch 4K Ultra HD TV (2016 Model)
 174



Cheetah Mounts APTMM2B TV Wall Mount for 20-75-inch TVs Bundle with 10-feet Braided...
 14,193



Sony HTXT2 2.1 Channel Sound Base with Bluetooth
 39
\$148.00

YouTube

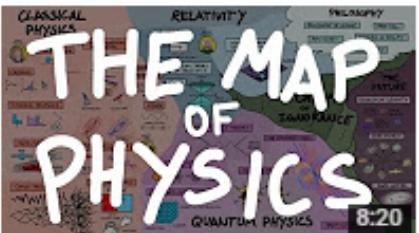
Recommended



[How the blockchain will radically transform the...](#)
TED 
126,953 views • 2 months ago



[How to gain control of your free time | Laura Vanderkam](#)
TED 
404,362 views • 3 weeks ago



[The Map of Physics](#)
DominicWalliman
604,613 views • 3 months ago



[สามก๊ก 2010 ตอนที่ 22 \(28 กุมภาพันธ์ 2560\)](#)
สามก๊ก 2010 TH
165 views • 2 hours ago



[The Future of Data Science - Data Science @ Stanford](#)
Stanford
48,516 views • 1 year ago



[จักรพรรดิมังโกลผู้พิชิตโลก จากยุคเริ่มต้นถึงล้มสลาย](#)
SFG Unicorn
38,487 views • 7 months ago



[10 อันดับสุดยอดขุนศึกในยุคสามก๊ก by CHERRYMAN](#)
CHERRYMAN
96,747 views • 3 weeks ago



[Android App dev by Visual Studio](#)
Suppakit Thongdee
5,859 views • 1 year ago

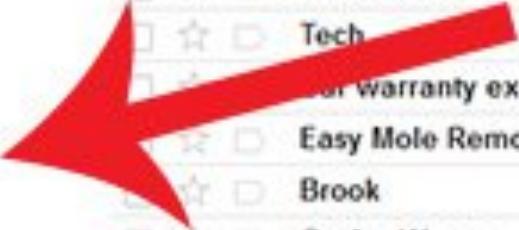
Gmail

Gmail •

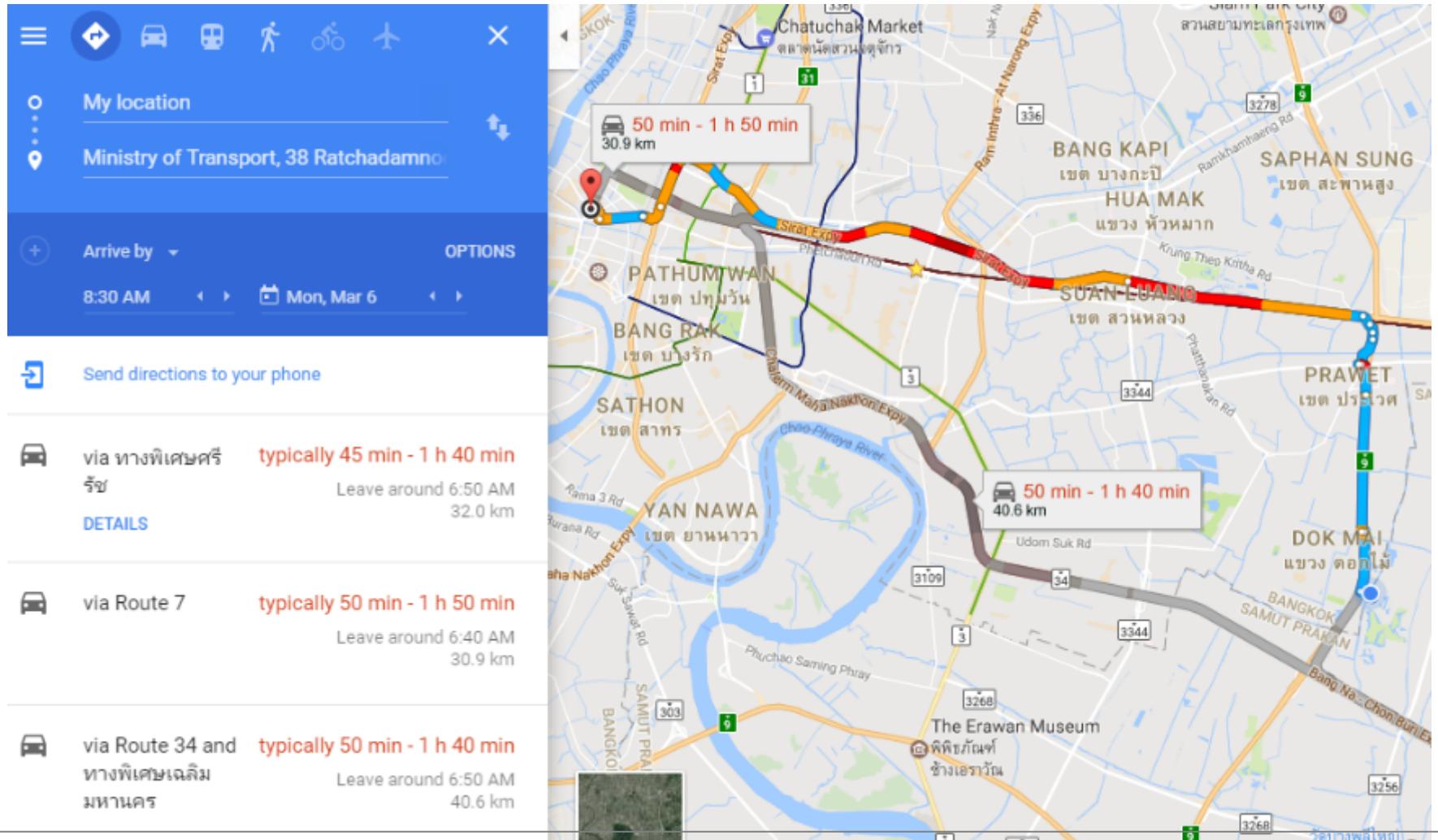
COMPOSE

Delete all spam messages now (messages that have been in Spam more than 30 days will be automatically deleted)

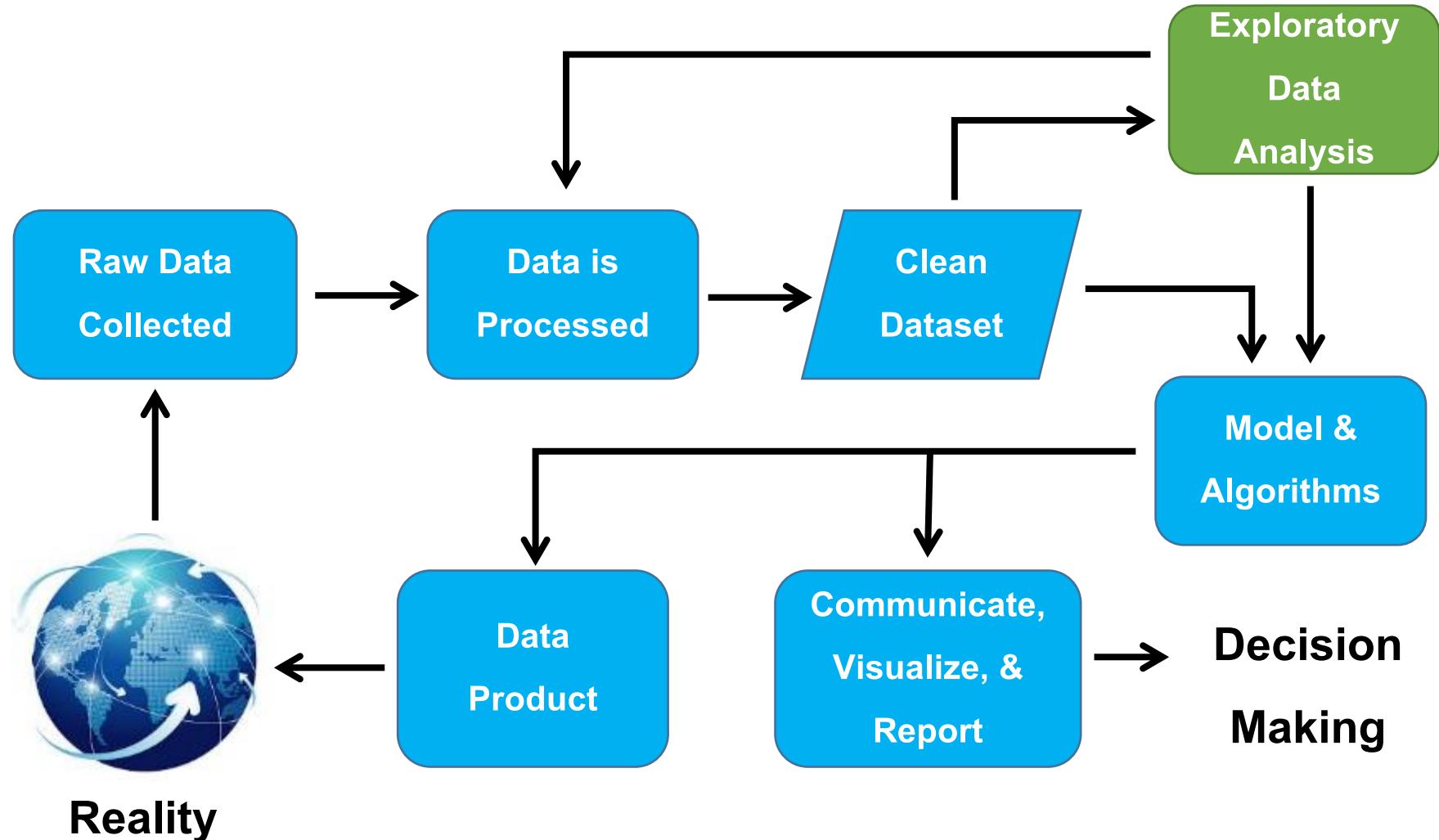
<input type="checkbox"/>	【家出少女を救う神待ち掲示板】	galen@ozdachs.com	家出少女を救う神待ちサ
<input type="checkbox"/>	BetterThanHCG (2)	galen@ozdachs.com	Traci says "It's BETTER"
<input type="checkbox"/>	Tech	galen@ozdachs.com	System Update - Click t
<input type="checkbox"/>	or warranty experts	galen@ozdachs.com	60% OFF - If you would li
<input type="checkbox"/>	Easy Mole Removal	galen@ozdachs.com	Remove Moles and Ski
<input type="checkbox"/>	Brook	Vmax Pills Official Site - 100% Guaranteed	4:23 pm
<input type="checkbox"/>	GetAnyWoman	galen@ozdachs.com	I got a date this weeke
<input type="checkbox"/>	Easy Mole Removal	galen@ozdachs.com	Remove Moles and Ski
<input type="checkbox"/>	LOTOTOjim	galen@ozdachs.com	••今月最後です••••□]
<input type="checkbox"/>	iPads Under One Hundred	galen@ozdachs.com	Absolutely, positively th
<input type="checkbox"/>	Jessica Iwane	galen@ozdachs.com	28 days later this 51 ye
<input type="checkbox"/>	Painting Services	galen@ozdachs.com	House need painting? I
<input type="checkbox"/>	Jessica Iwane	galen@ozdachs.com	HAVE YOU SEEN THIS:
<input type="checkbox"/>	Cobra Health	galen@ozdachs.com	Cobra Health for galen

 A large red arrow points from the left towards the second email in the list, specifically highlighting the subject line "Tech".

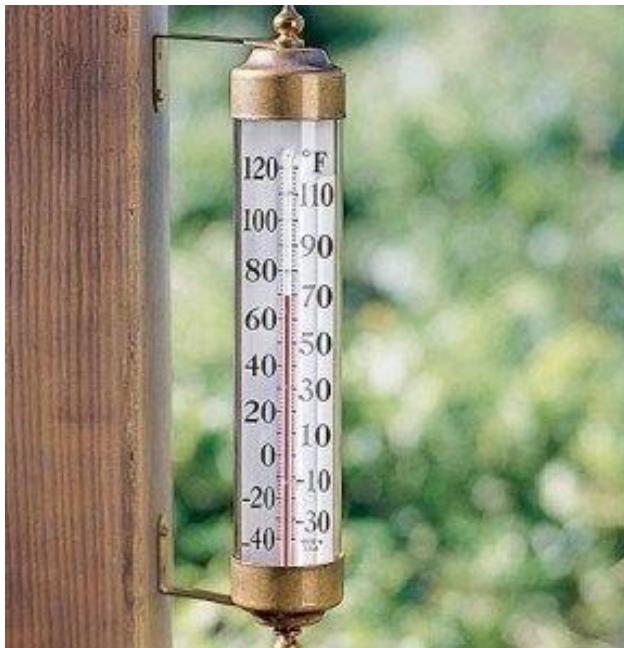
Google Map



Data Science Process



Reality



Raw Data Collected

ମୁଖ୍ୟ ପରିକାଳିକା



#	Time	App / Site (1 to 10)	Category	Tags	bulk edit
1	8h 40m	Dreamweaver	Dev Tools	webdev work	[edit]
2	4h 37m	mail.google.com/a/	Comm (Email)	all-comm google-apps work	[edit]
3	3h 31m	Photoshop	Design/Presentation	design webdev work	[edit]
4	2h 36m	mail.google.com	Comm (Email)	all-comm google-apps work	[edit]
5	1h 30m	news.ycombinator.com	News/Blogs	personal	[edit]
6	1h 23m	twitter.com	Social Networking	personal social	[edit]
7	1h 10m	localhost:3000	Dev Tools	webdev work	[edit]
8	45m 20s	rescuetime.com	Personal Productivity	webdev work	[edit]
9	36m 3s	google.com/header	News/Blogs	google-apps work	[edit]
10	34m 58s	deck.rescuetime.com	Dev Tools	design webdev work	[edit]

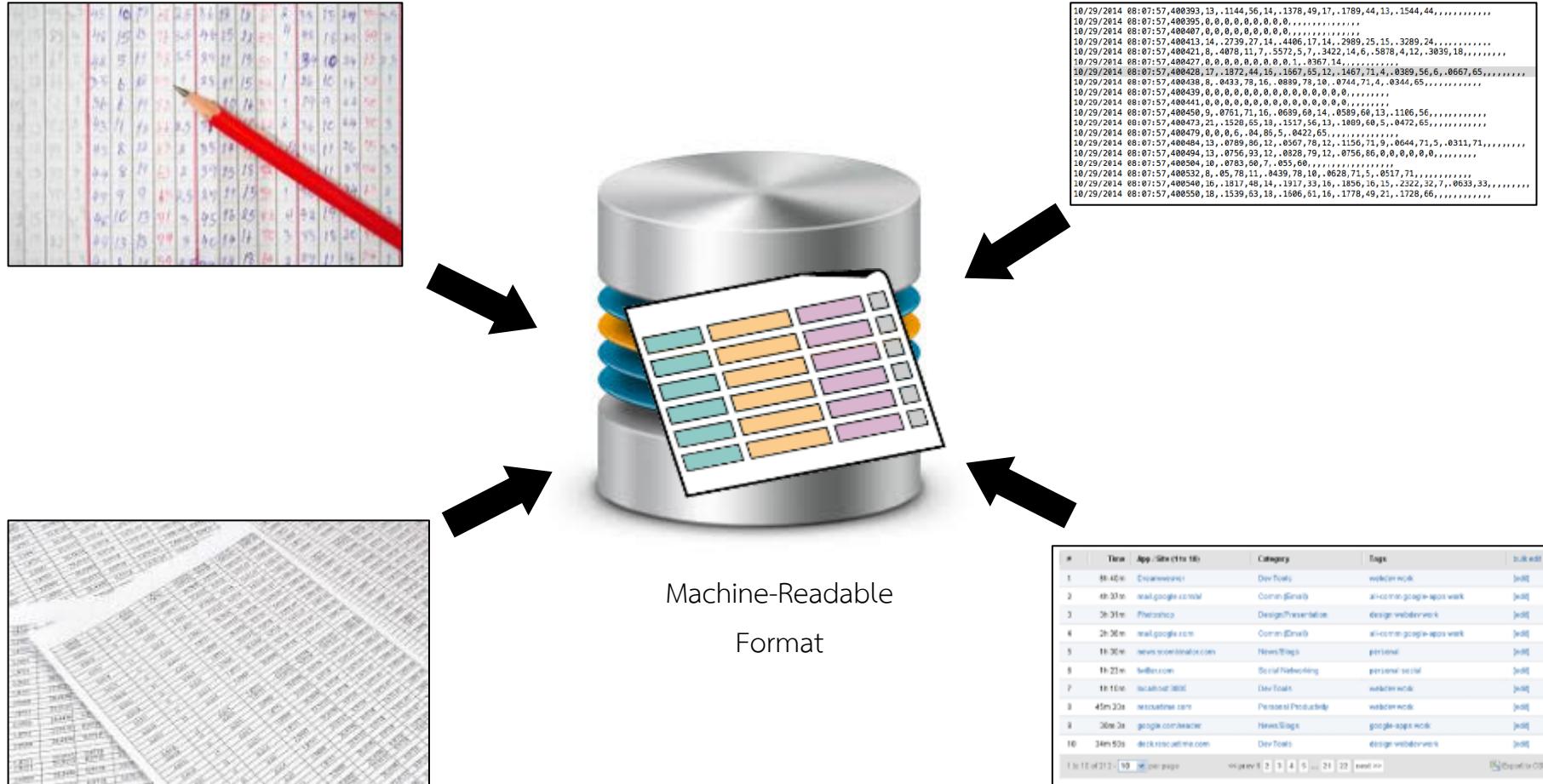
1 to 10 of 212 - 10 per page

[<< prev](#) [1](#) [2](#) [3](#) [4](#) [5](#) ... [21](#) [22](#) [next >>](#)

 Export to CSV

Data is Processed

ମୋବାଇଲ୍ ସିଗ୍ନାଲ୍ ଫର୍ମେଟ କୁ computer ରେ ଆପଣଙ୍କିରି



- Ref:** • (content) URL
• (image) URL

Data is Processed

- Merge Data Sets into the same format
 - Rebuild Missing Data with appropriate values
 - Standardize *ສຳເນົາຂອງຕົວຢ່າງ* e.g. same column name
 - Normalize *ກົດລູກຄະຫຼາດ* e.g. same date format
 - De-Duplicated
 - Verify & Enrich e.g. update the salary values
- (ມີຕົວຢ່າງ)*

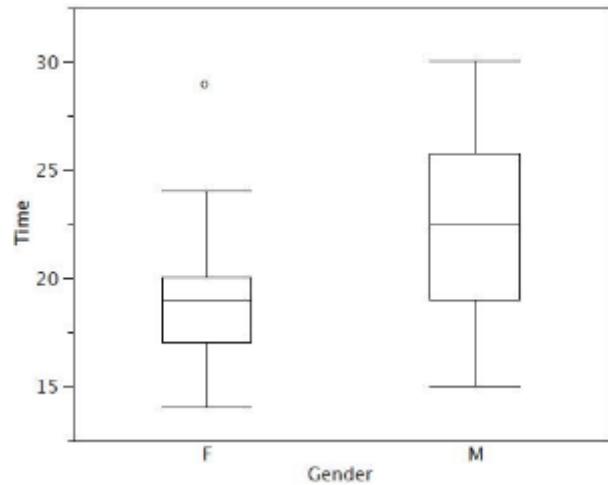
Clean Dataset

ໄດ້ອະນຸຍາກສູງ clean dataset

- Good Format
- Good Shape
- Ready for Data Analysis

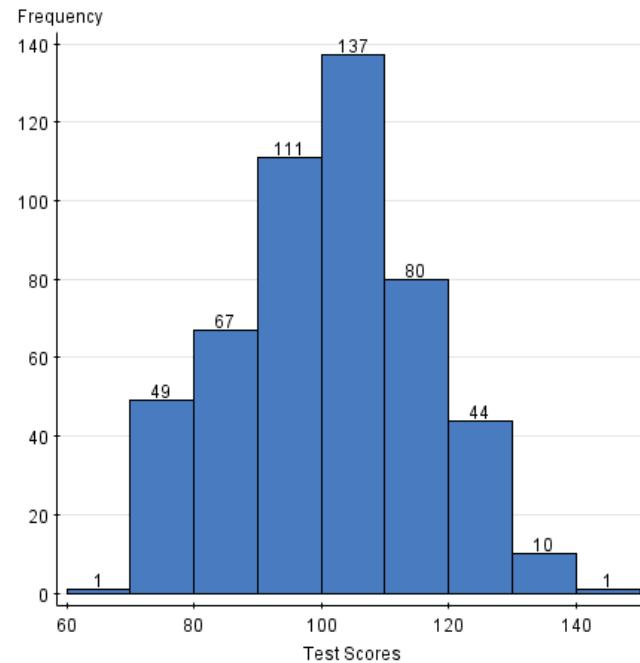
Exploratory Data Analysis

ជំរួចរាល់



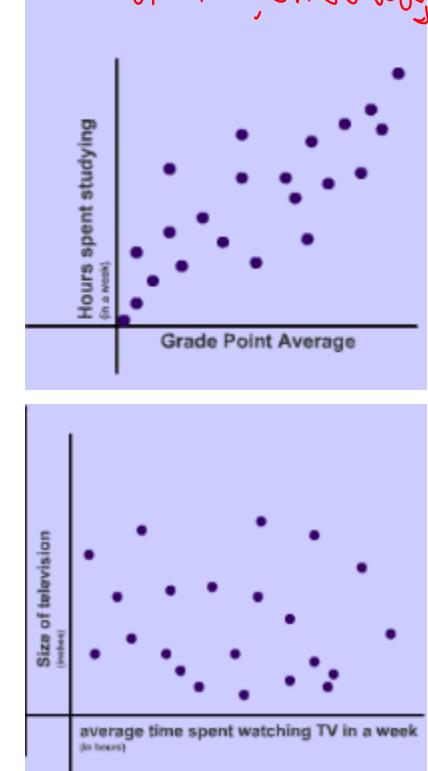
Box Plot

គម្រោងនៃការបកចែកសីវិទី
នៅក្នុងភេទប្រឈាន



Histogram

គម្រោងសរុប
នៅក្នុងតម្លៃដែលគម្រោង



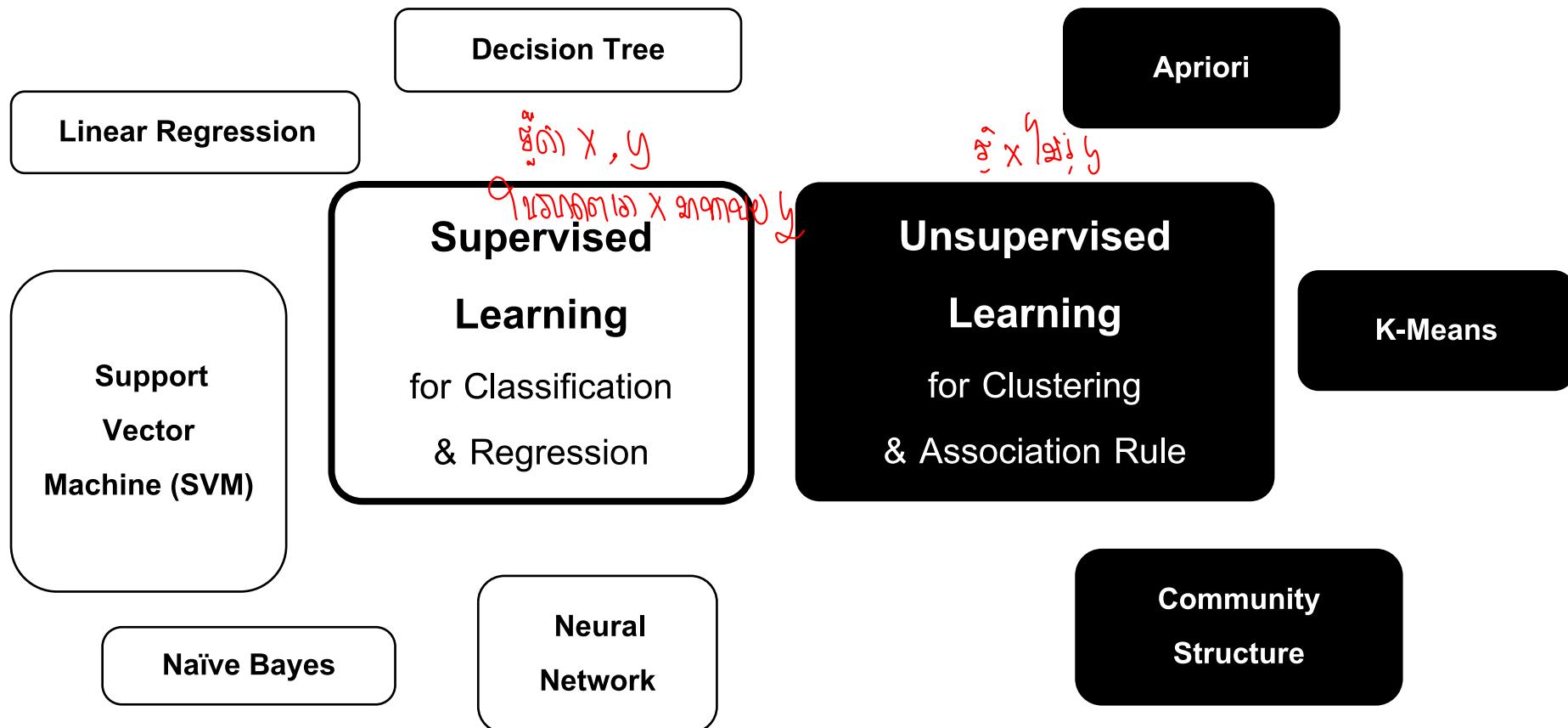
Scatter plot

គម្រោងសម្រាប់
ជូន

Model & Algorithms

- **Algorithm** → ក្រែបាយនៃរាជការណ៍ដែលមិនត្រូវគុណភាពខ្លះ → tree class 9 is 10Tree ? នៅ
- a process or set of rules to be followed in calculations or other problem-solving operations, especially by a computer.
- **Model** → សំណងសម្រាប់គុណភាពដែលត្រូវបានបង្កើតឡើង → មិនមែនបច្ចុប្បន្ន
- A model is a computation or a formula formed as a result of an algorithm that takes some values as input and produces some value as output.
- **Example**
 - Model: Decision Tree with structure
 - Algorithm: A process to build an accurate decision tree

Model & Algorithms (Example)

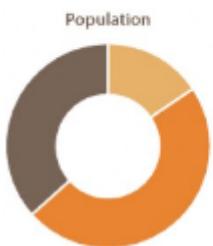
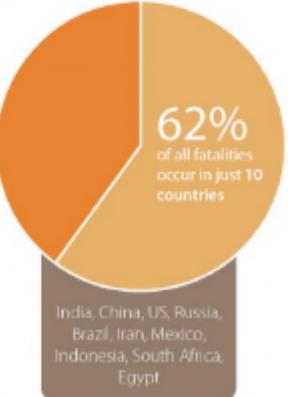
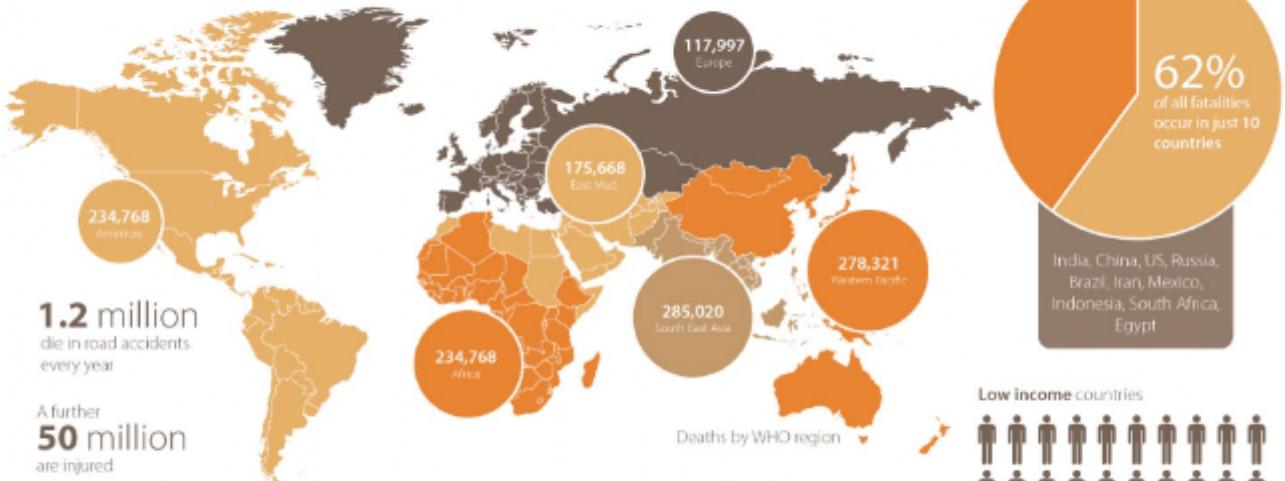


Communicate, Visualize, & Report

(រៀនវិ)

Road Traffic Accidents: The Modern Killer

The Global Status Report released by WHO this year, confirms that road traffic injuries are still a big global health and development problem



- High income countries
- Middle income countries
- Low income countries

90% of deaths occur in low or medium income countries

However high income countries have over **50%** of all registered vehicles



The Laws



Only **49%** of countries stipulate a **legal blood alcohol concentration** limit of less than 0.05g per decilitre



Only **57%** of countries **requires seatbelts** to be used by passengers



Only **40%** of countries have a comprehensive **helmet law** and require helmets to be of a specific standard

On the rise?

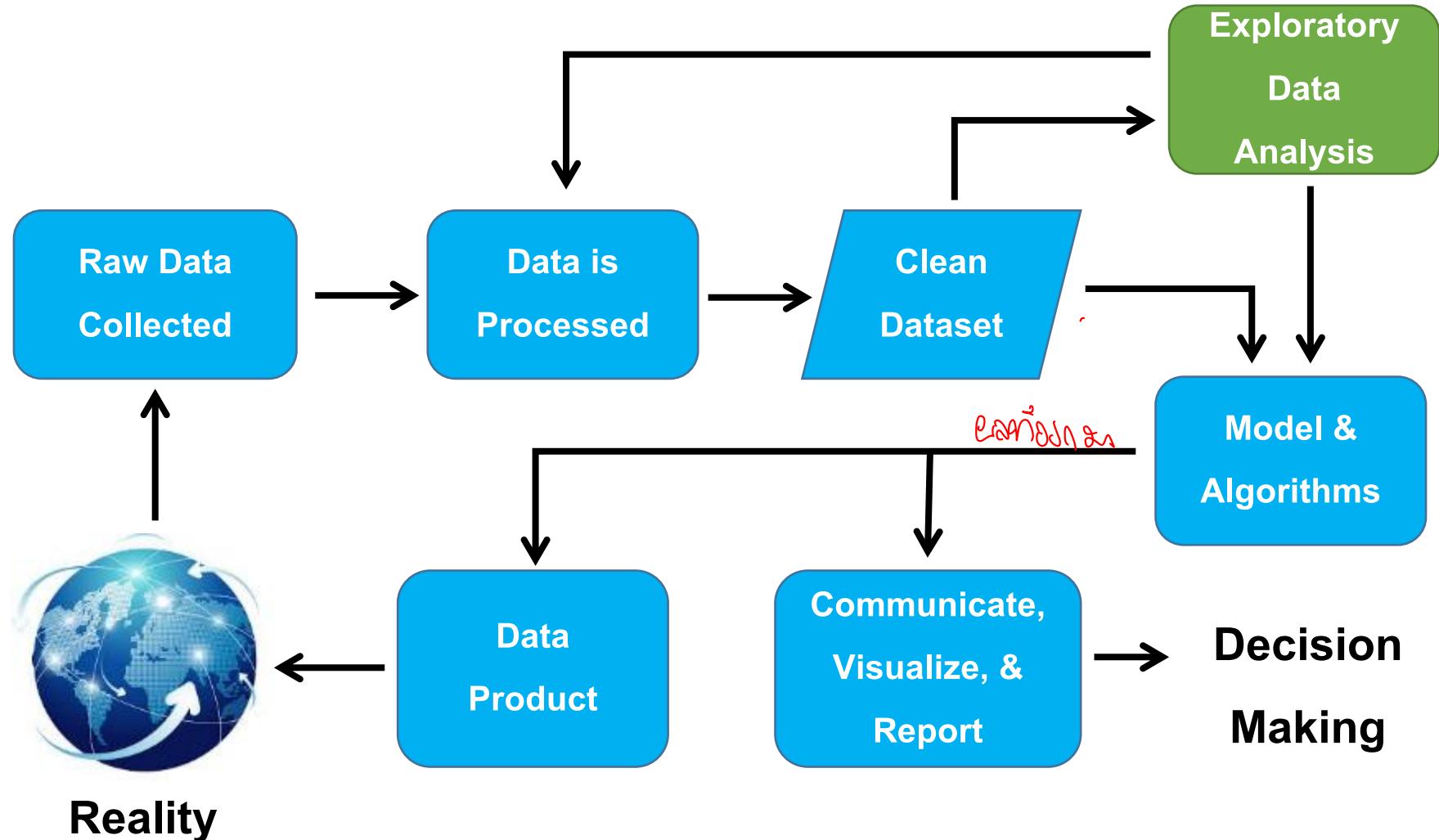
Road traffic accidents are predicted to rise to the **5th leading cause of death by 2030**, higher than AIDS, lung cancer and diabetes



Car accidents are the ...
number 1 killer for 15-29 year-olds

Data Product

Data Science Process



MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants

DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

PROGRAMMING & DATABASE

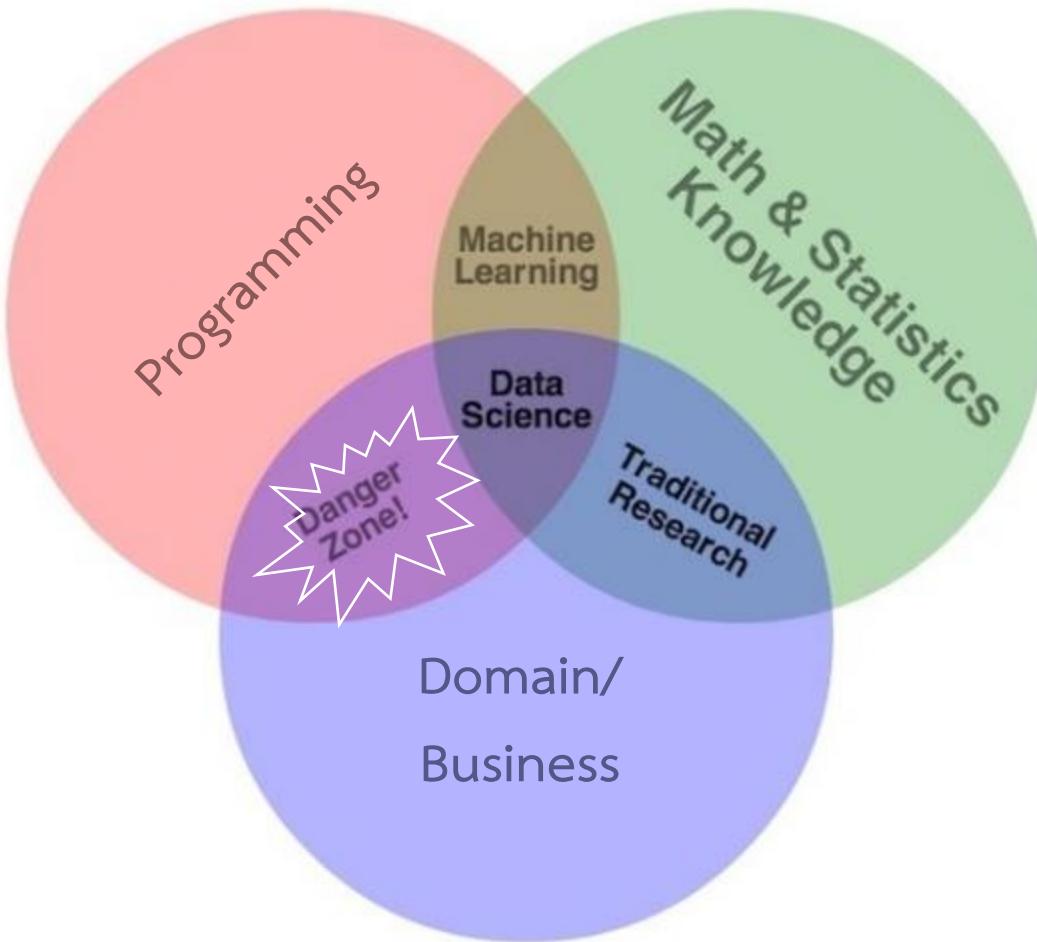
- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS



COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

One Last Thing



“

The goal is
to turn data into information,
and information into insight.

”

Carly Fiorina

つづく