

CPE 213 Data Models

(a.k.a. Data Modeling and Visualization)

Lecture 1: Introduction to Data Modeling and Visualization

Asst. Prof. Dr. Santitham Prom-on

Department of Computer Engineering, Faculty of Engineering
King Mongkut's University of Technology Thonburi

Course Learning Outcome

- Evaluate and apply suitable data modeling techniques to analyze real-world data.
- Create meaningful visualization that address the relevant problems.
- Understand the data science process and the role of data scientists.



Reference Textbook

- Grolemund, G. and Wickham, H. (2017). R for Data Science. O'Reilly Media. Link: <http://r4ds.had.co.nz>
- Slides, papers, and additional documents will be provided in-class and online

Grading

- Midterm Examination 30%
- Final Examination 30%
- Project ~~20%~~ 15%
- Lab / HW (10 LABS) ~~20%~~ 10%
- Quest 15%

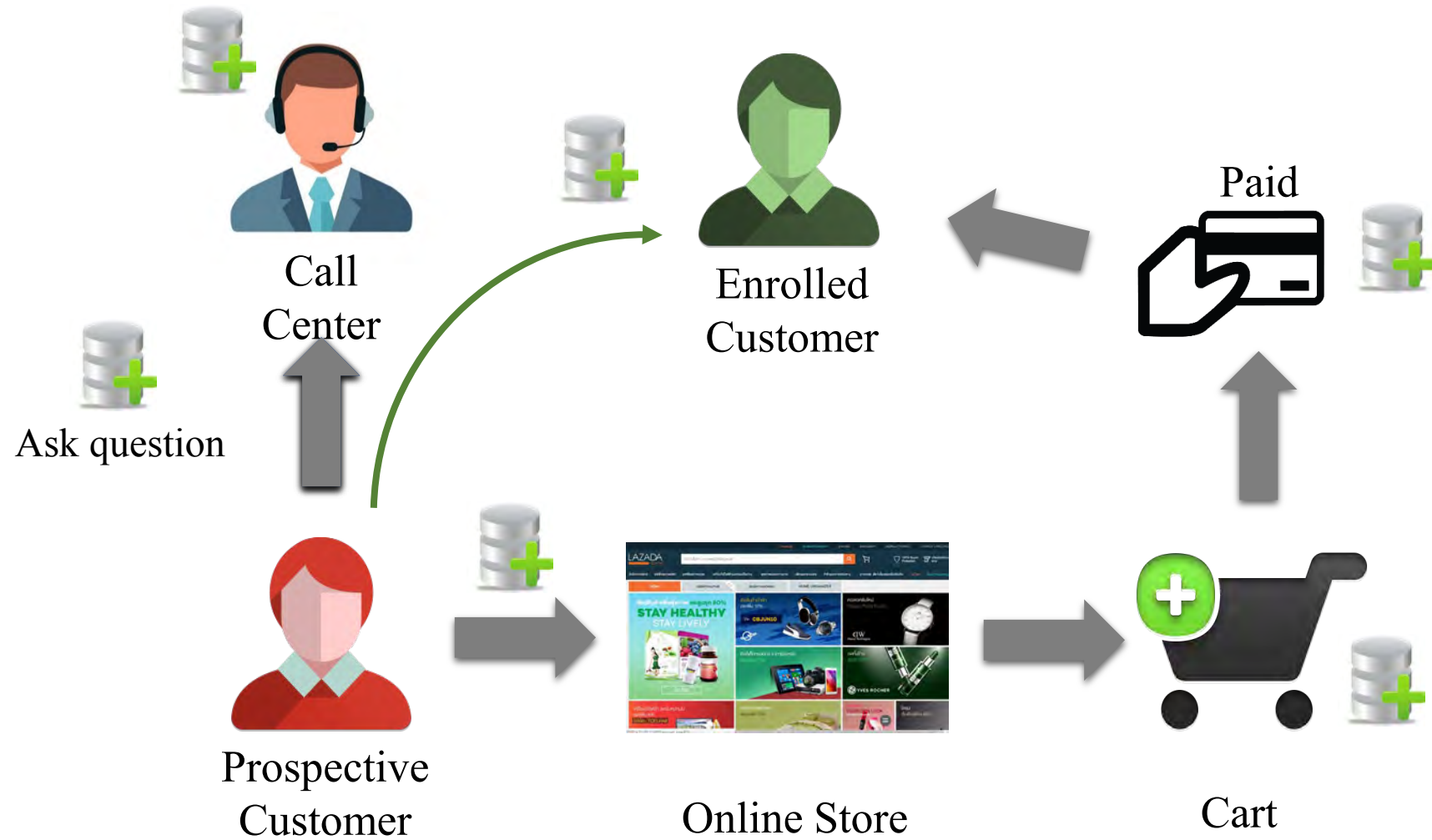
#	Date	Topic
1	22 Jan 2021	Introduction to data modeling and visualization
2	29 Jan 2021	R programming
3	5 Feb 2021	Basic types of data visualizations
4	19 Feb 2021	Data visualization 1 - distributions
5	5 Mar 2021	Data visualization 2 - relationship
8		=== Midterm Examination ===
9	19 Mar 2021	Data visualization 3 - network
10	26 Mar 2021	R database operations
11	2 Apr 2021	Modeling statistical distribution
12	9 April 2021	Linear regression
13	23 April 2020	Logistic regression
14	30 April 2020	Decision tree
15	7 May 2020	Term project presentation
16		=== Final Examination ===

Project

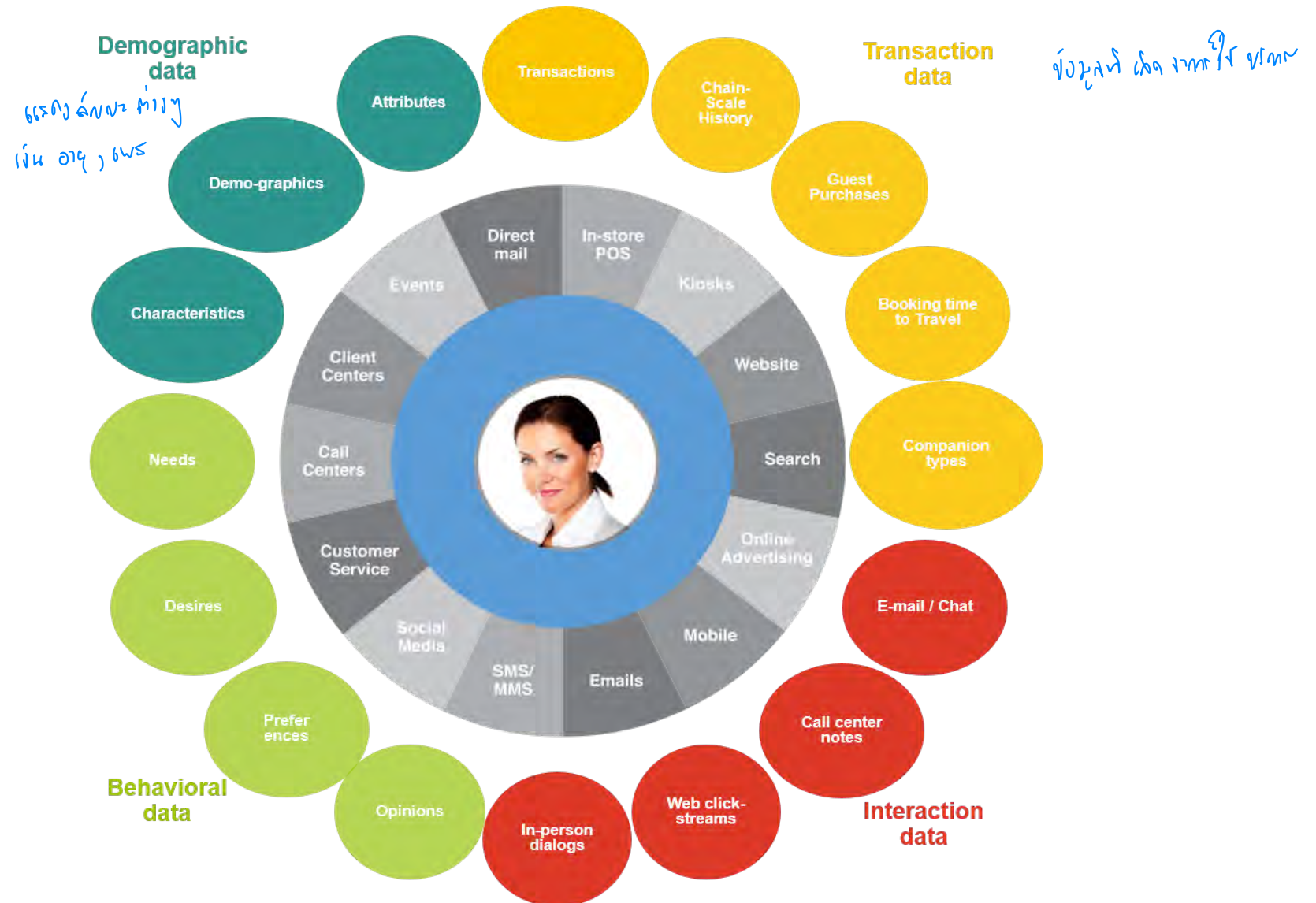
- Students work in teams of 5
- Task: Visualization
 - Acquire a dataset and understand it
 - Create a visualization of a dataset
 - Address related questions
- Task: Modeling
 - Identify the modeling objective
 - Build a model with the prepared data
 - Make predictions based on the model
- Output: Write a technical report and present the visualization and model to the class
- Due: 7 May 2021

What are Data Modeling and Visualization?

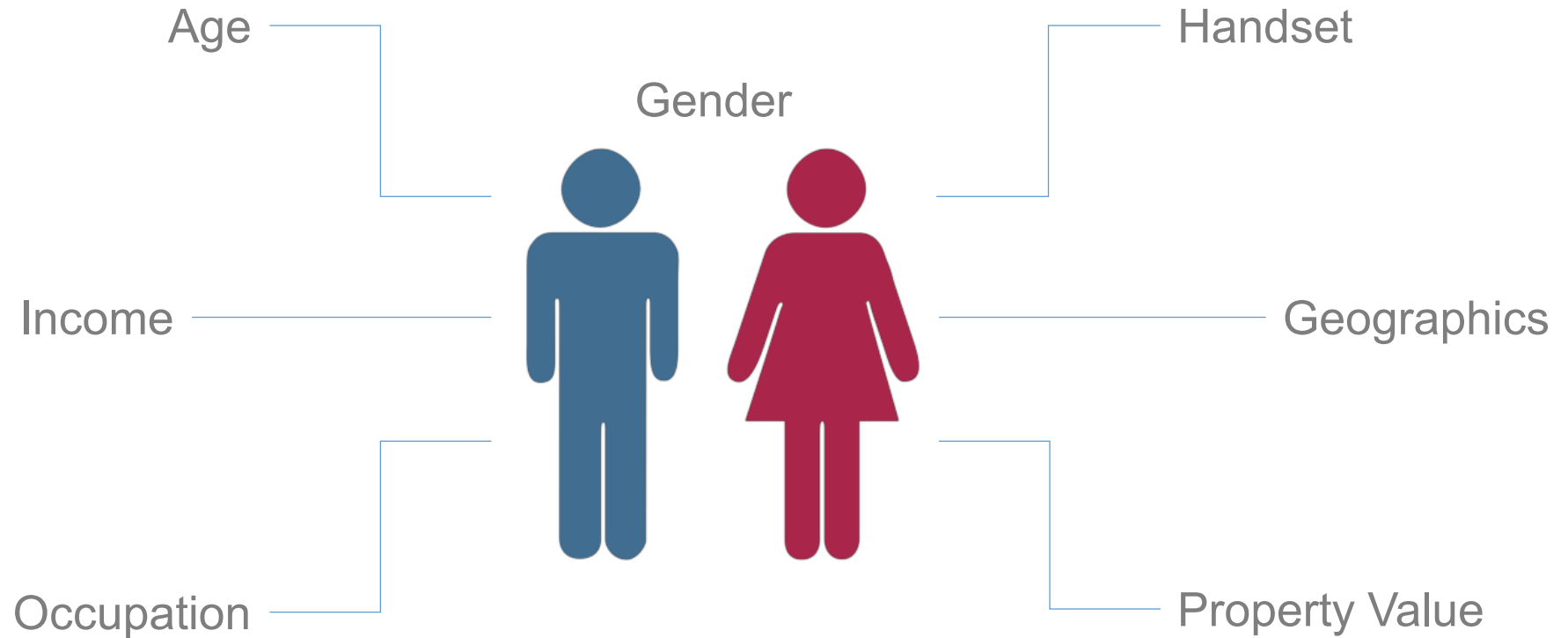
Section



CUSTOMER JOURNEY



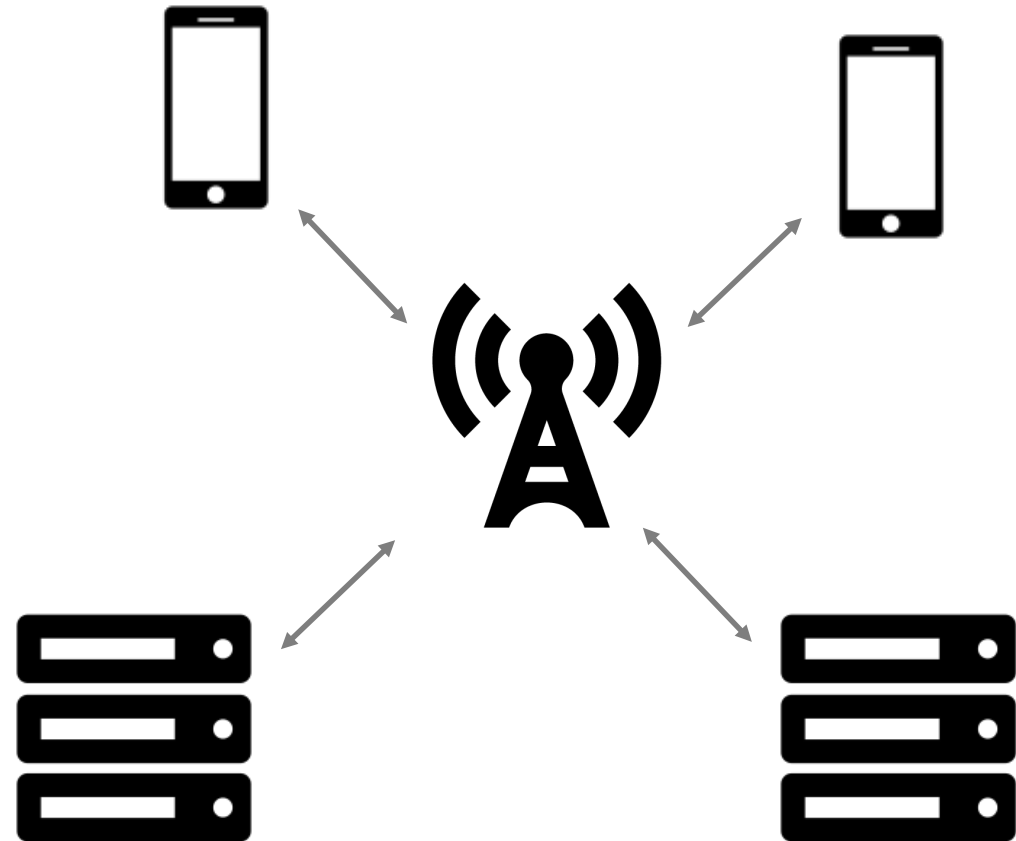
Demographic data



Transactional data

(information service)

- Call detail record
- Internet usage
- Tower location
- CRM data
- Service plan
- Payment information
- App usage
- Call center



අනුකූල වෙනම තත්ත්වයන් ගබඩා කිරීමේ අවස්ථාව

Interaction data

- E-mail log
- Chat log
- Web click stream
- App click stream
- In-person log



Visitor Jane

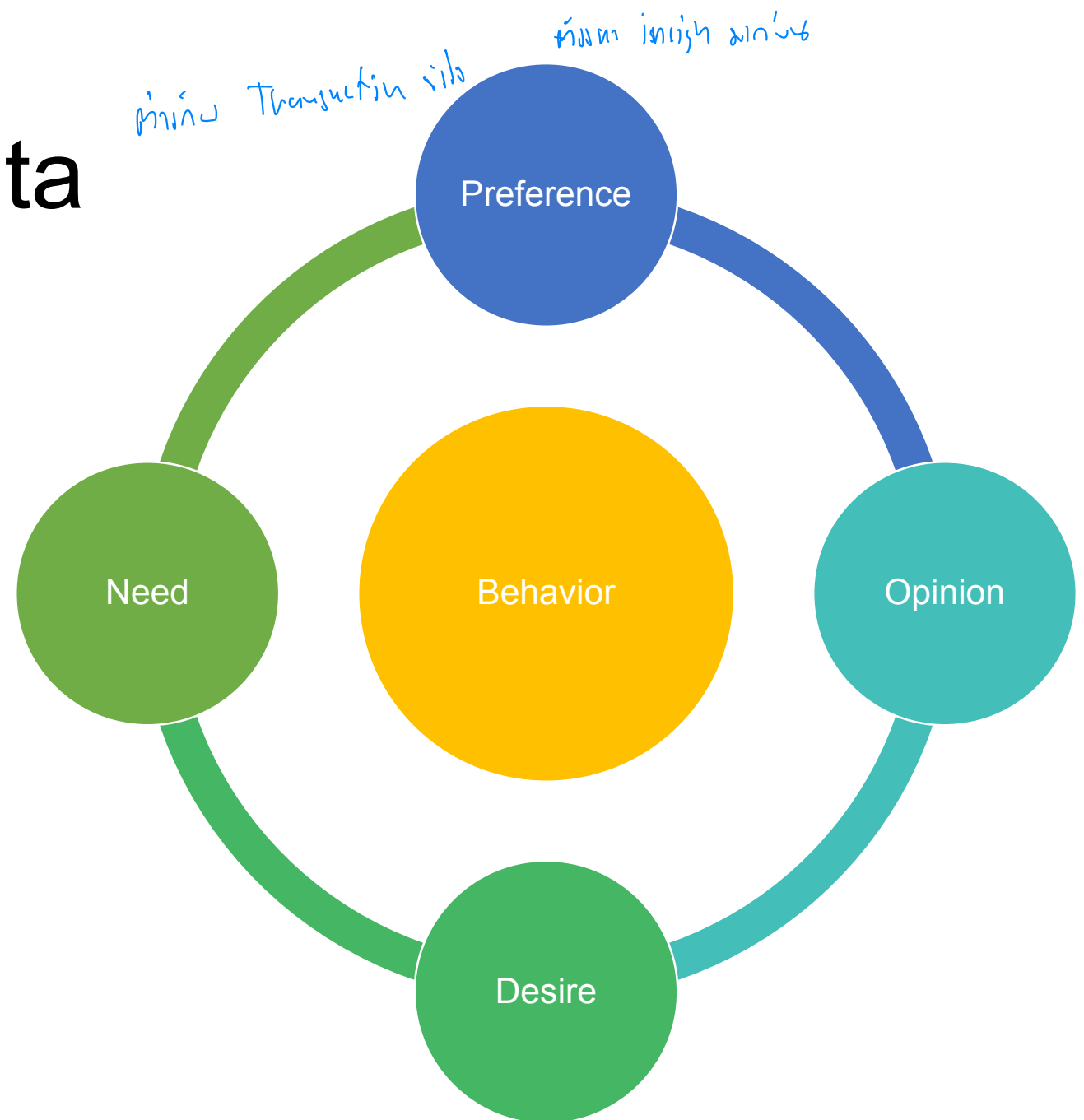
Hello. Do you offer free shipping?

Hi Jane, yes, we offer
free local shipping!



Behavioral data

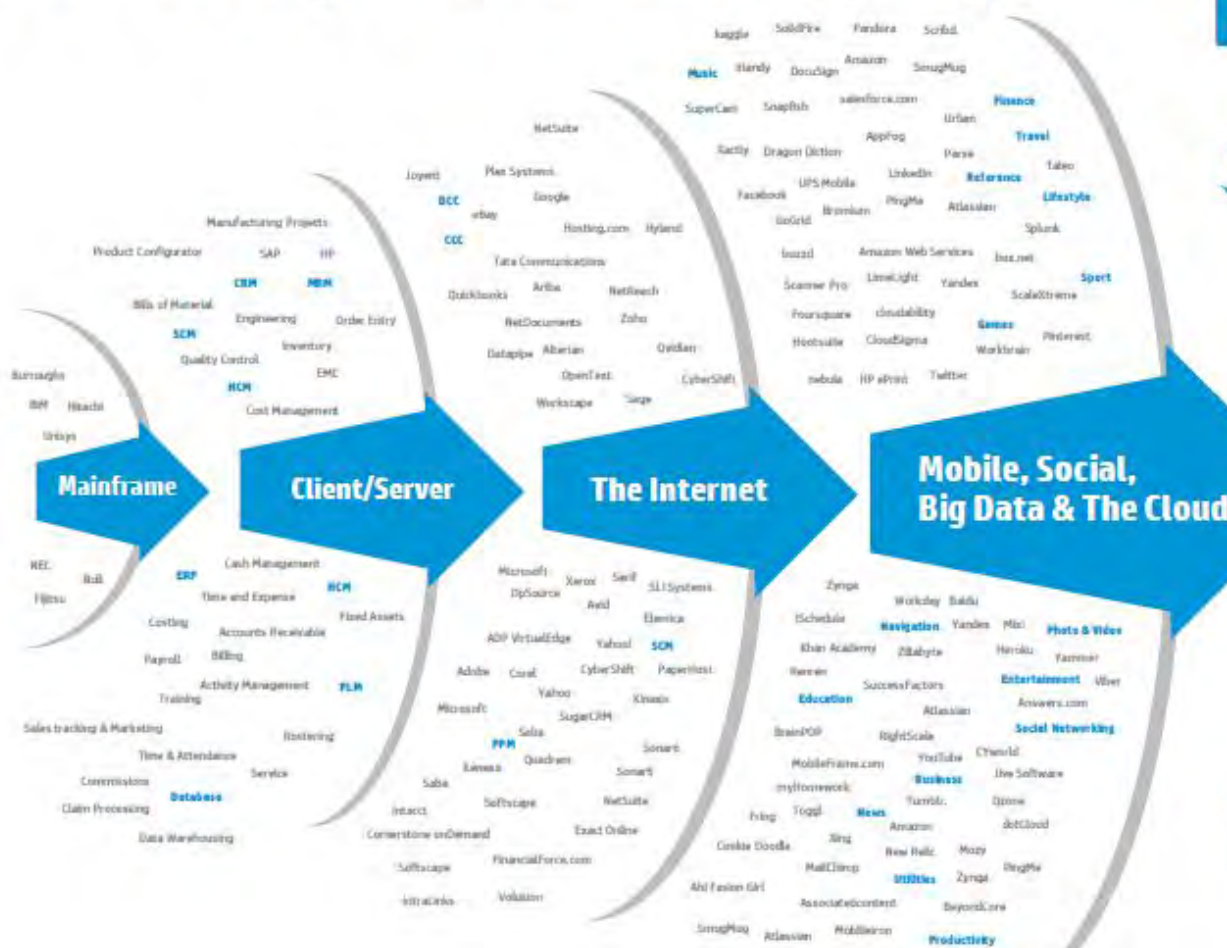
- Opinion
 - Like matrix
- Preference
 - Collaborative filtering
- Desire
 - Wish list
 - Browsing history
- Need
 - Tax/Law/Regulation
 - Life
 - Finance





The rise of data ...

A new style of IT emerging



Every 60 seconds



98,000+ tweets



695,000 status updates



11 million instant messages



698,445 Google searches



168 million+ emails sent



1,820TB of data created



217 new mobile web users

"DATA IS THE NEW OIL."

From the beginning of recorded time until 2003, we created **5 exabytes** of data.

In 2011 the same amount was created **every two days**.

By 2013, it's expected that the time will shrink to **10 minutes**.

Every hour, we create enough Internet traffic to fill **7 billion DVDs**.

Side by side, that's **seven times** the height of **Everest**.

Coined in 2006 by Clive Humby, a British data commercialization entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considered data to be an economic asset, like oil.

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

As of August 2012, there were just over

4 million articles in the English Wikipedia.

There are **133 million BLOGS** on the web.

80%

of all humans own a mobile phone of some sort. Out of 5 billion mobiles, 1 billion are smartphones. (In Singapore, 54% of citizens are smartphone users.)

English is the dominant language of the web. But by 2014 it will be **Chinese**, if its current rate of increase continues.

Top languages used on the web (May 2011):



247 billion EMAILS

are sent **every day**. (Up to 80% are spam.)

10% of all photos ever taken were taken in 2011.

60%

of all humans (5.4 billion people) are active texters. In 2010, 193,000 text messages were sent **every second**.

high-frequency traders

Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan, with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

These specialized algorithms make split-second decisions to buy or sell a commodity. New cable being laid under the Atlantic will shave

5 milliseconds

from the current 65 milliseconds it takes for trading instructions to travel between New York City and London.

With new fiber-optic cable,

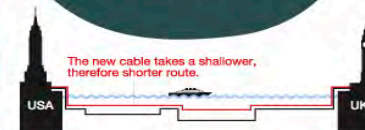
the round-trip time between New York and London will be 59.6 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cable (and who will pay millions to do so).

How they save 5 milliseconds

The depth of the Atlantic Ocean varies.

The new cable will lie on areas of the ocean floor that are up to 1,000 feet shallower than the current fastest cable. By taking a different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.



50% of 5-year-old kids in the U.S. are given access to a smartphone.

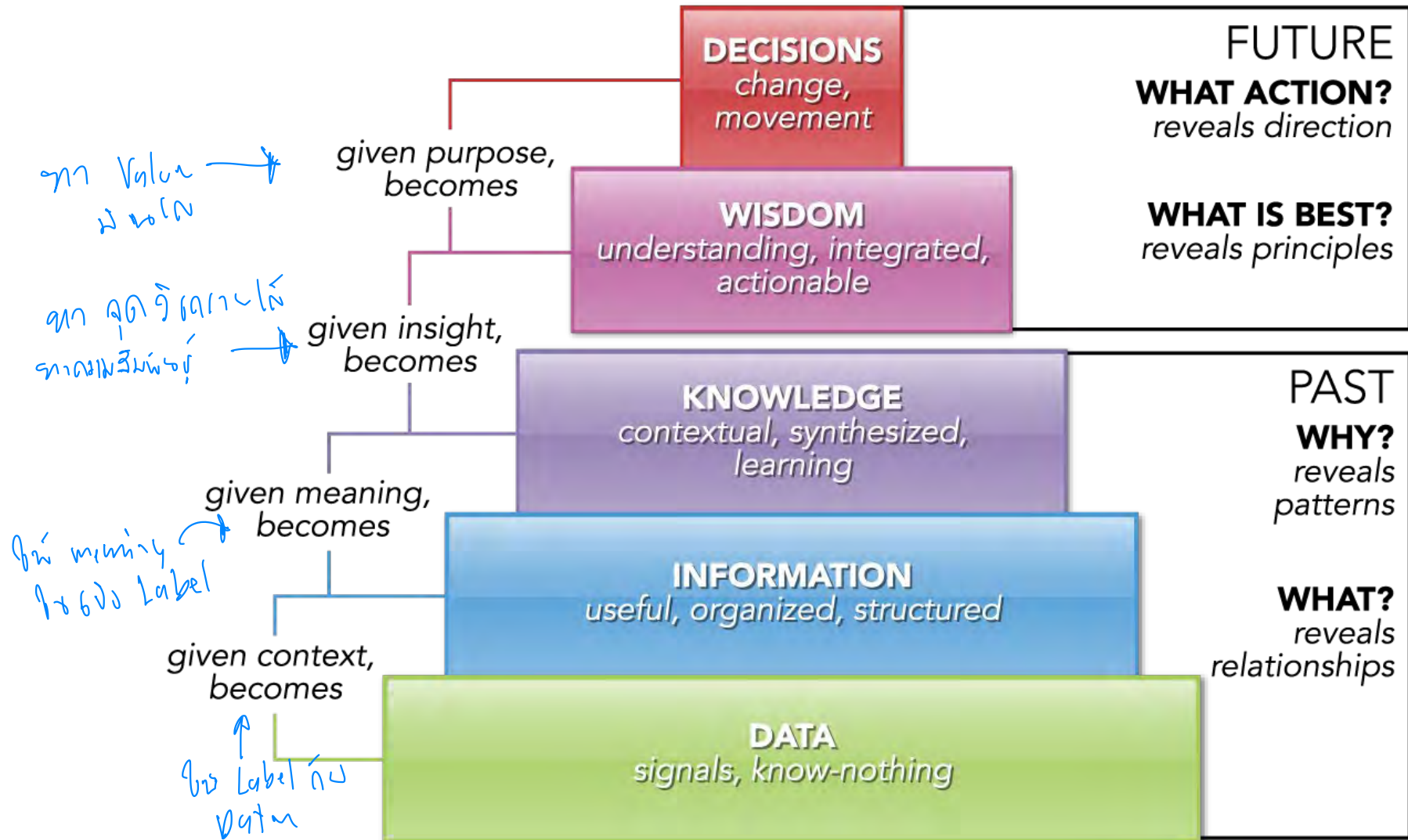


DATA

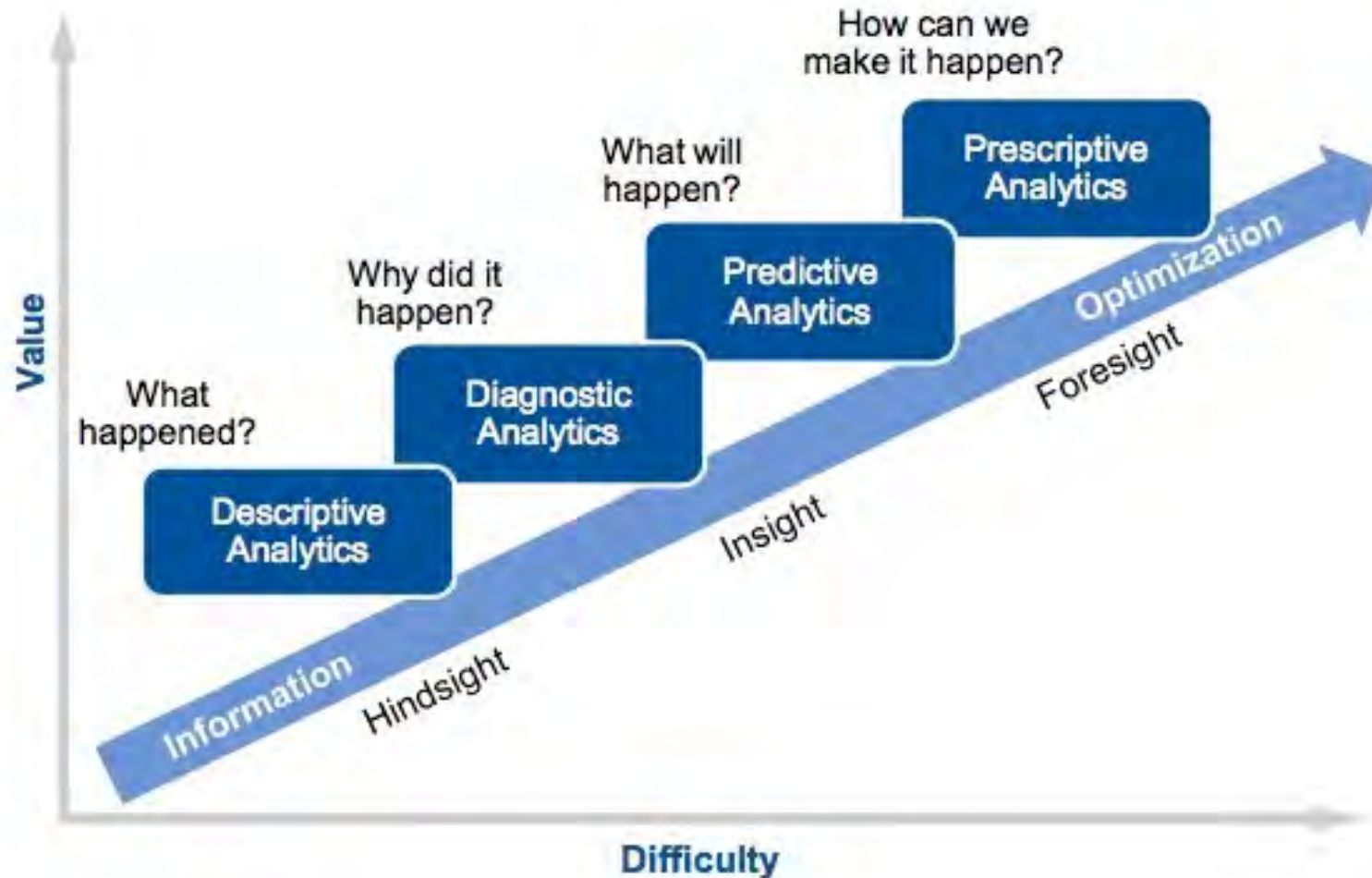


DECISION

DIKW (D) Pyramid

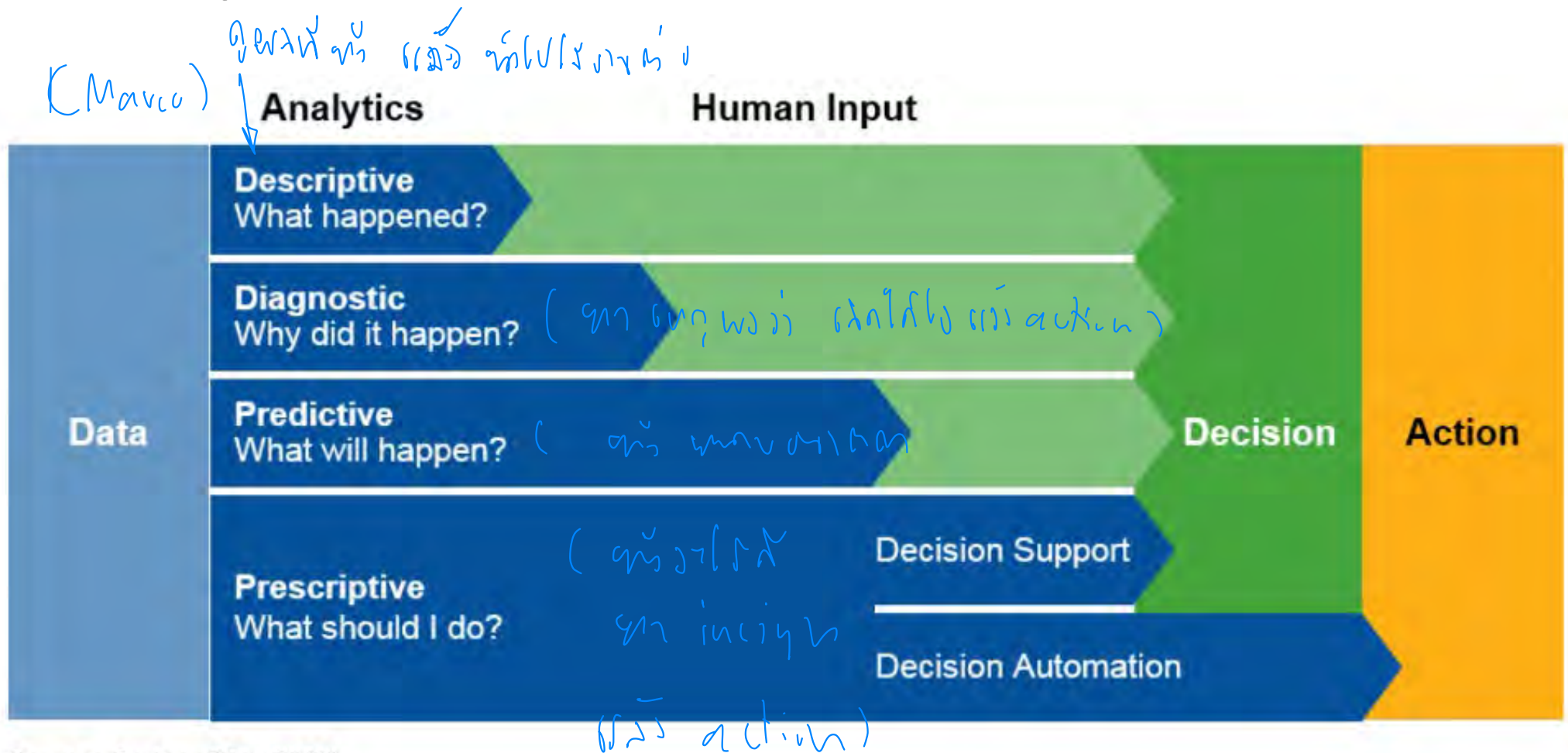


From descriptive ... to prescriptive



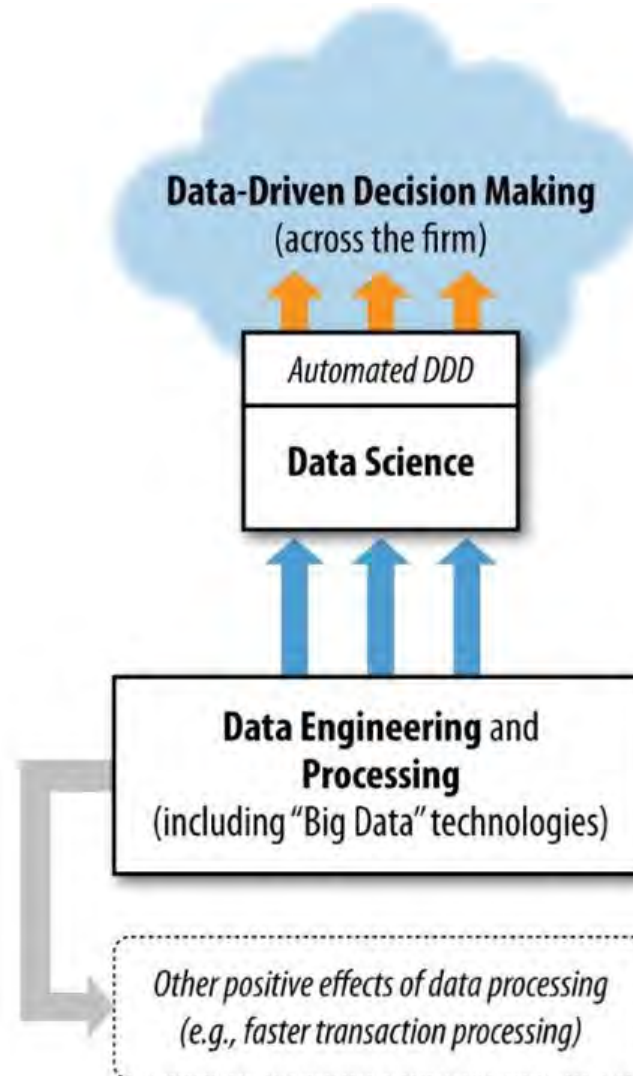
Source: Gartner (March 2012)

Analytics Capabilities Framework



Source: Gartner (May 2015)

Data-driven decision making



DDD = practice of basing decision on the analysis of data, rather than intuition

Principles and techniques for understanding phenomena via the analysis of data.

Accessing and processing of massive-scale data flexibly and efficiently with Big Data technologies

The Synopsis

The science

Extracting useful knowledge from data to solve business problems can be treated systematically by following a process with reasonably well-defined stages.

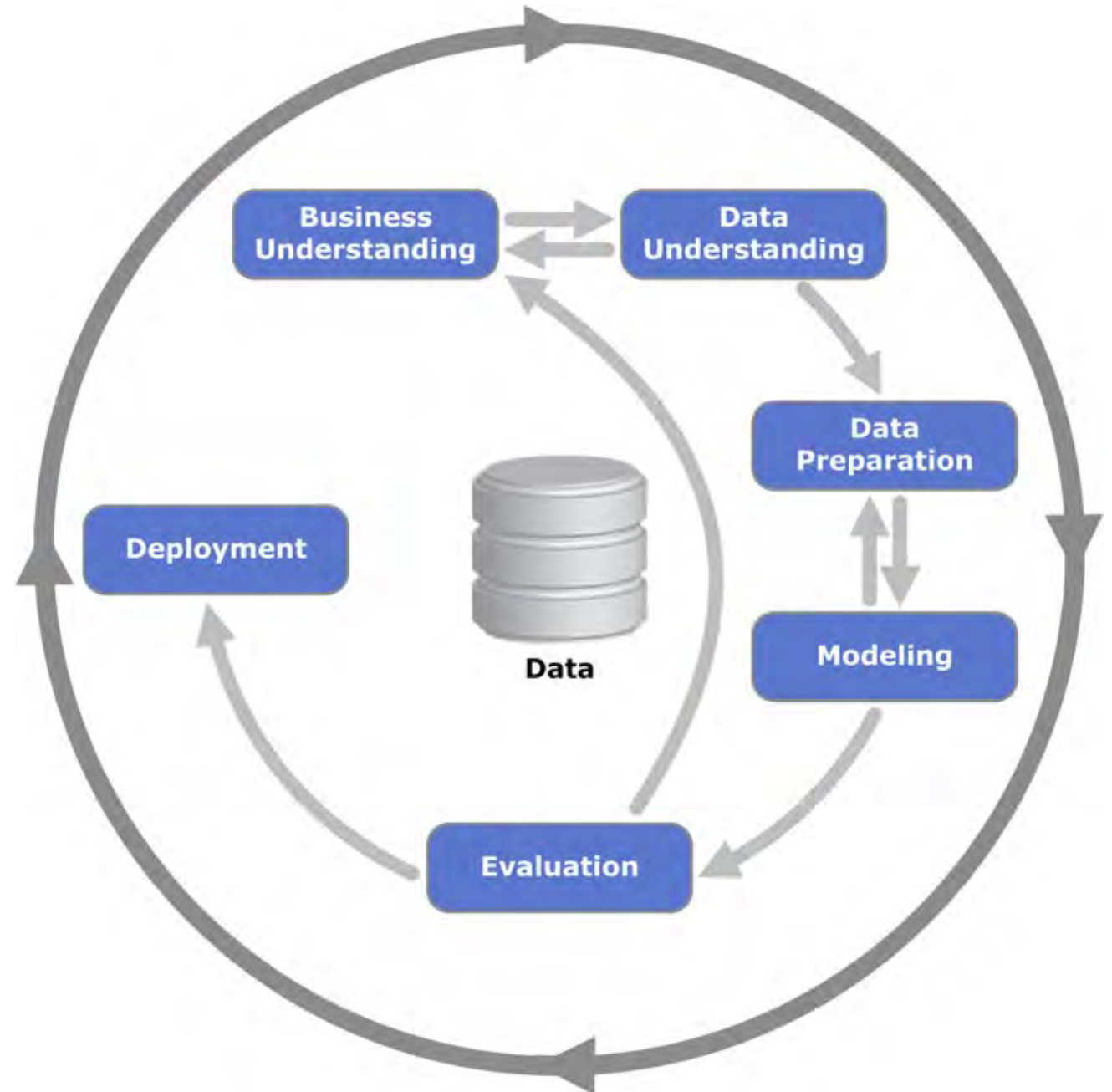
The technology

From a large mass of data, IT can be used to find informative descriptive attributes of entities of interest

Steps in data analytics

Cross-Industry Standard Process for Data Mining

CRISP - DM



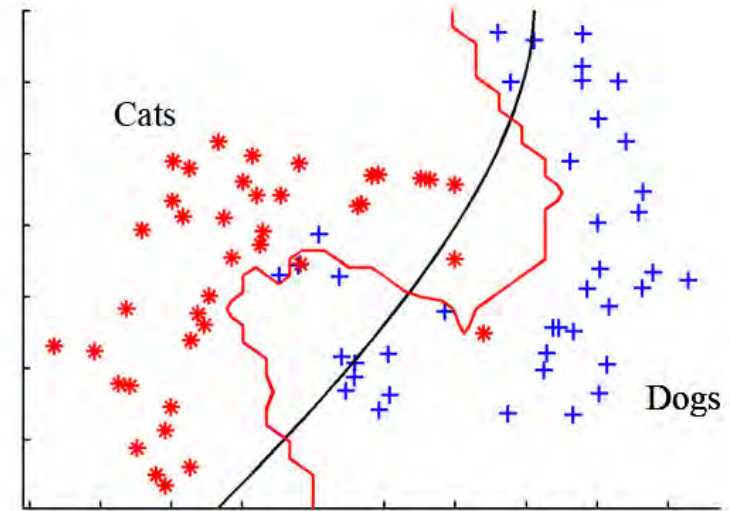
MACHINE LEARNING MODEL

Machine Learning

“The science of getting computers to learn from data without having to be explicitly programmed by humans.”

Machine learning is surrounding you

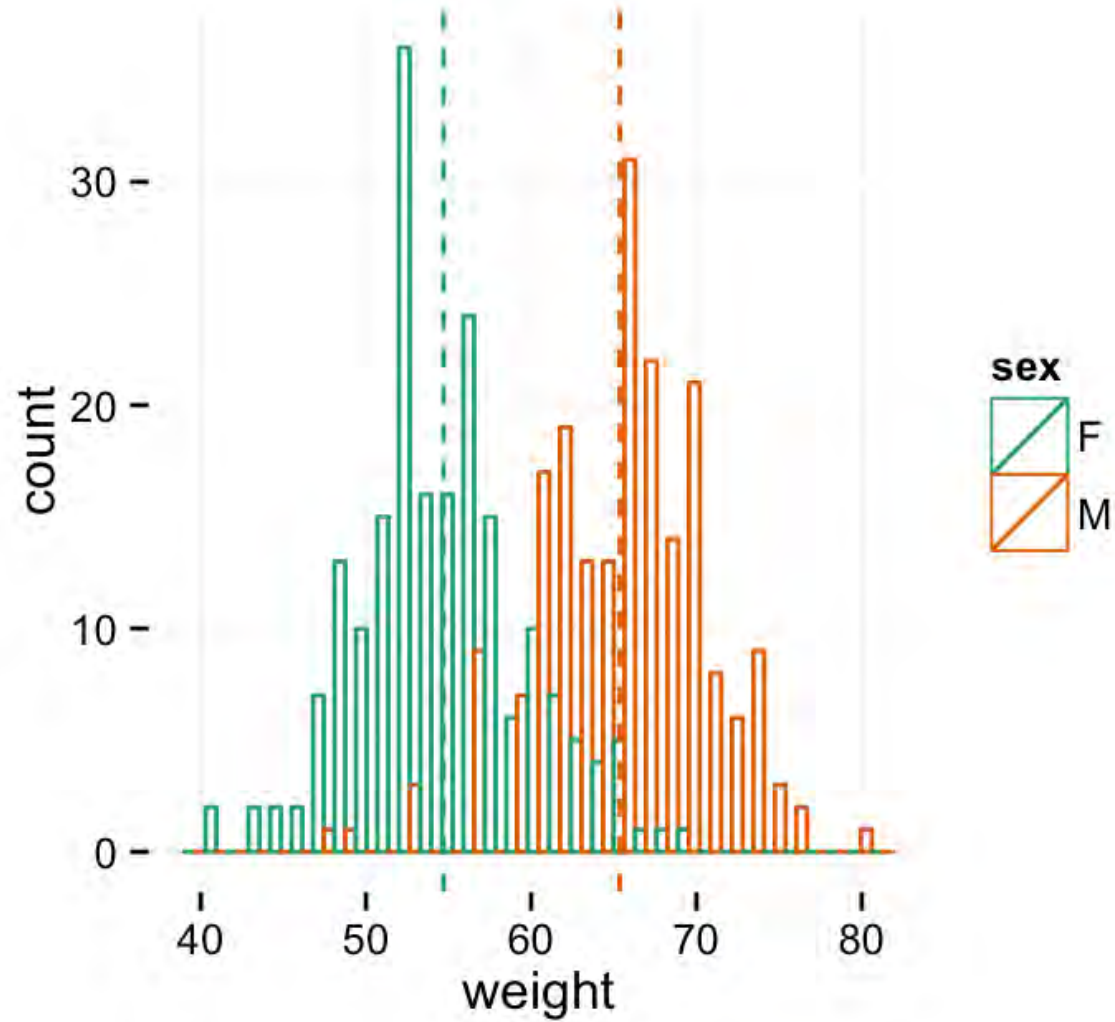
- Google search
- Auto Facebook photo tagging
- Email Spamming
- Games
- Chat bot
- Recommender



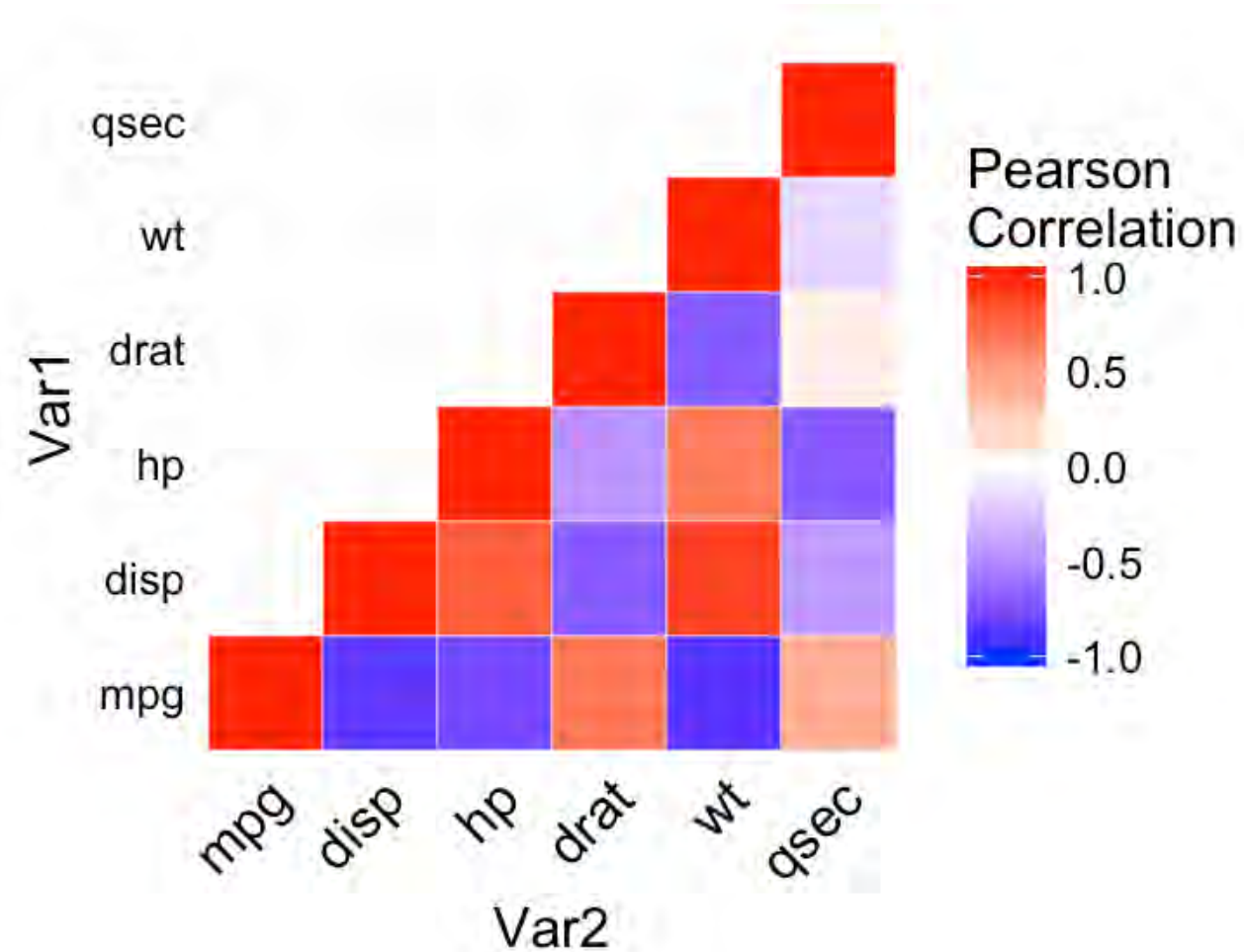
Visualization



Visualization

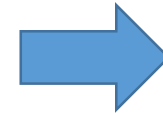


Visualization

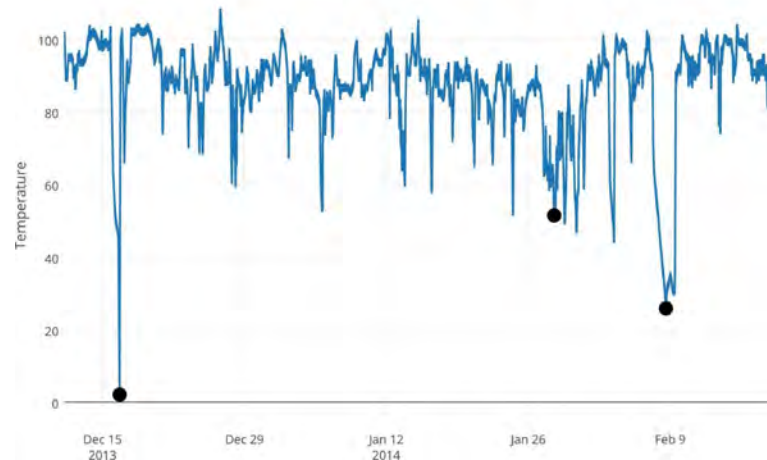


Use Case

Anomaly detection Problem



Products



Pain Point

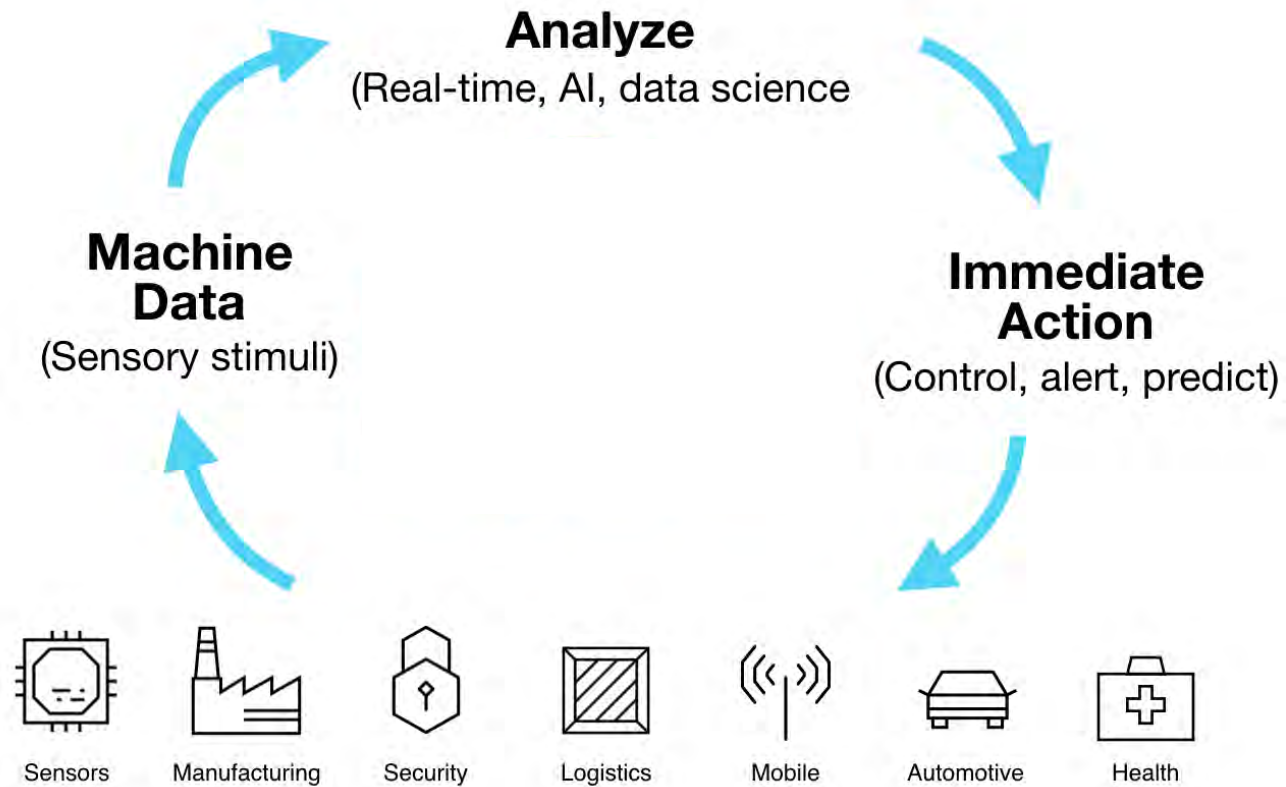
Sometime anomalies in production occur and result in fail QC

Anomaly detection

Analytic objective

To monitor manufacturing parameters and
detect the event that is not normal

Anomaly detection Machine data



Anomaly detection

Machine data

Time

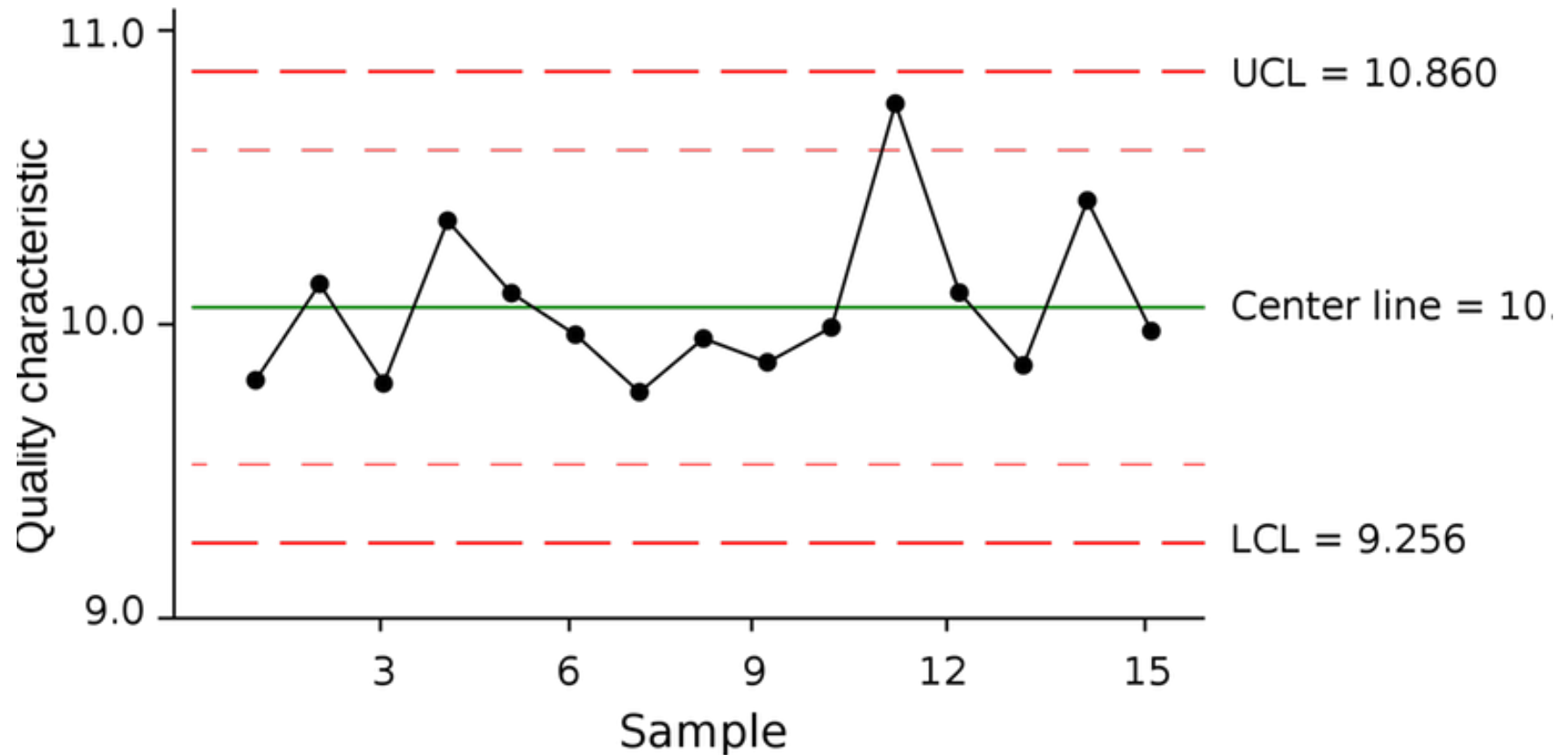


Sensor

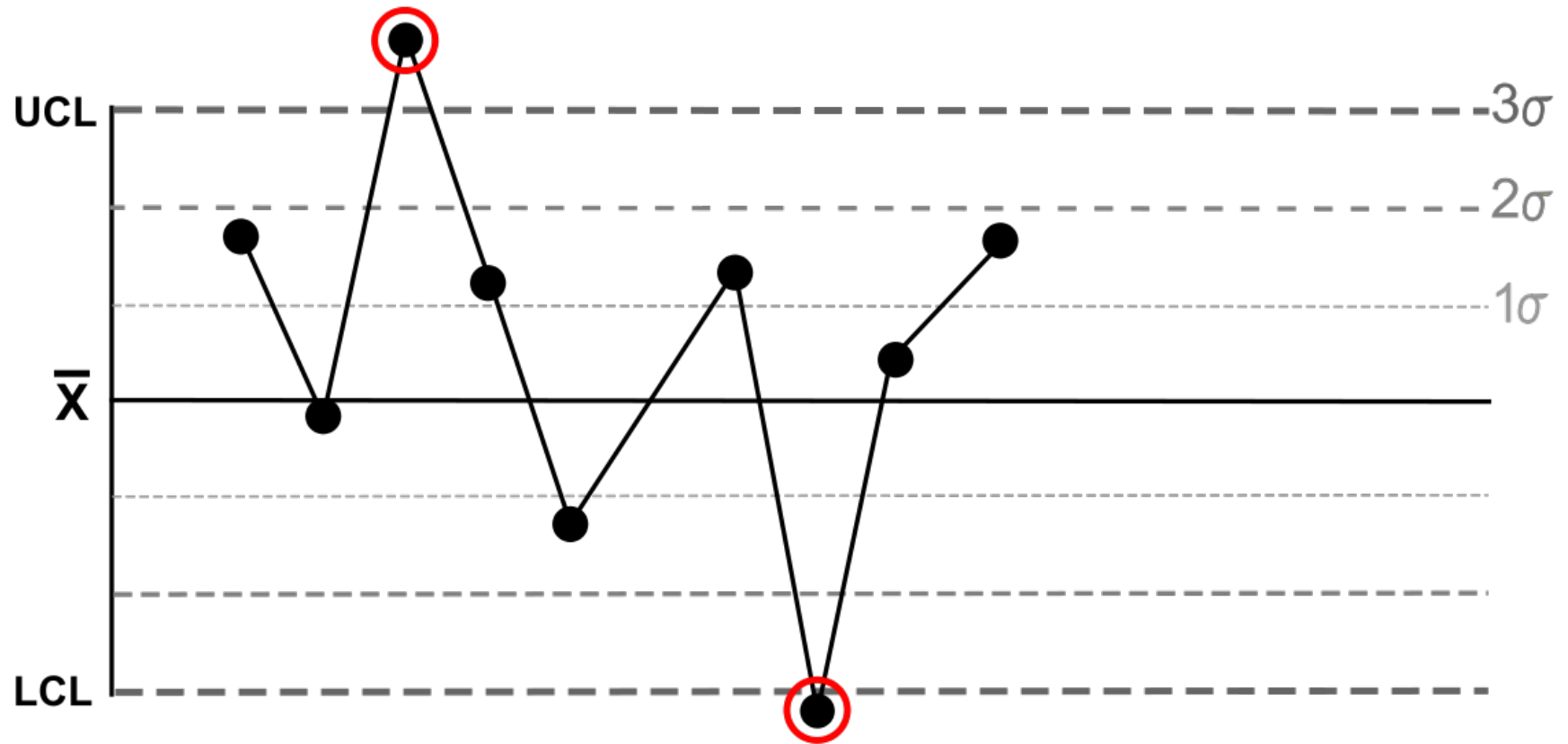
date	time	Sensor						output
		Inputv1	Inputv2	Inputv3	Inputv4	Inputv5	Inputv6	
8/29/2018	19:50:00	1	0	0	1	0	0	1
8/29/2018	19:55:00	1	0	0	1	0	0	0
8/29/2018	20:00:00	1	0	0	1	0	0	1
8/29/2018	20:05:00	1	1	1	0	0	0	1
8/29/2018	20:10:00	1	1	1	0	0	0	1
8/29/2018	20:15:00	1	1	0	1	0	0	1
8/29/2018	20:20:00	1	1	0	1	1	0	0
8/29/2018	20:25:00	1	0	0	1	1	0	1
8/29/2018	20:30:00	1	0	0	1	1	0	1
8/29/2018	20:35:00	0	0	0	1	0	0	1
8/29/2018	20:40:00	1	0	0	1	1	0	1
8/29/2018	20:45:00	0	0	0	1	0	0	0

Anomaly detection

Statistical process control



Anomaly detection



Anomaly detection

Outcome and usage

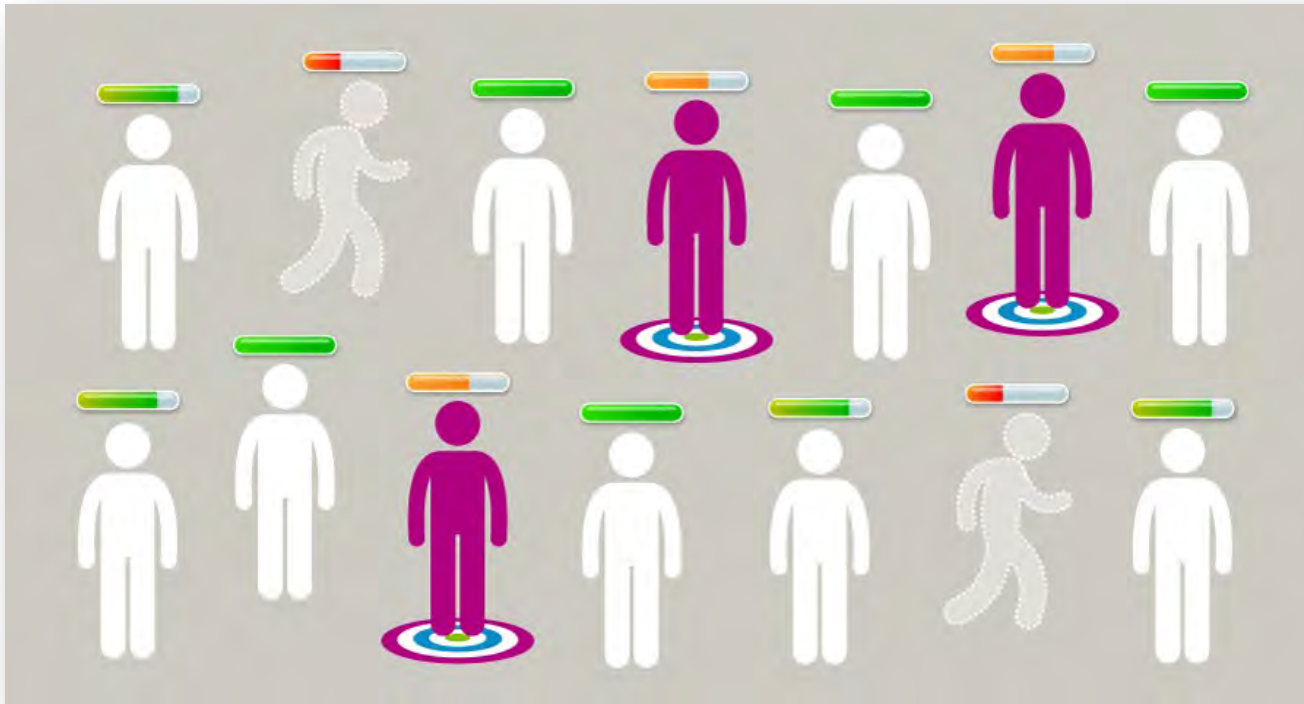
Usage

- Monitor the manufacturing parameters
- Detect and alert managers/operators once anomalous events are detected

Outcome

- Able to intervene with manufacturing process in time

Churn prediction Problem



Churn



Drop in Revenue

Churn prediction Analytic objective

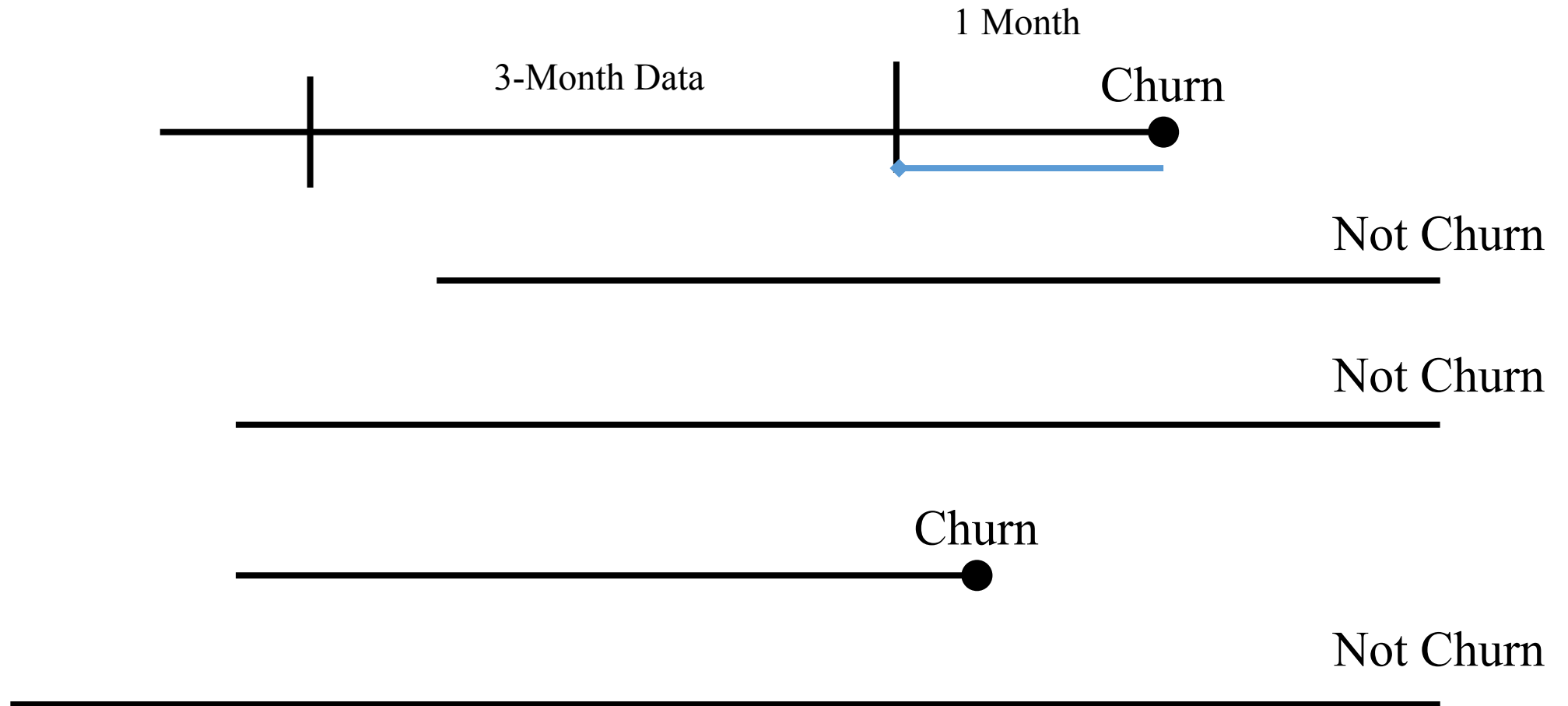
Who are likely to churn?

Churn prediction Timeline



PROBLEM: we know very little about customers

Churn prediction Data collection



Churn prediction Data types

Age
Service Month
Handset Type
Main Pack

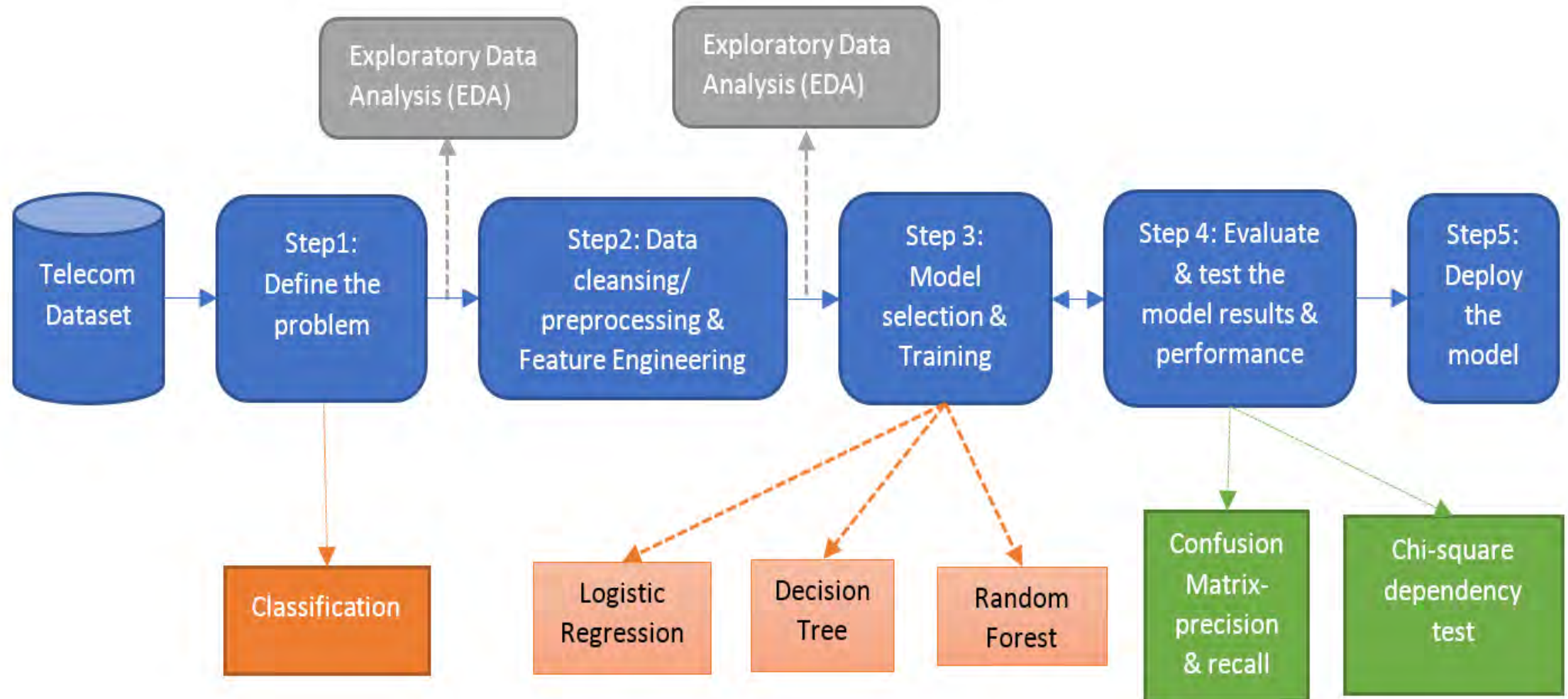
Product purchase
Browsing history
Preference



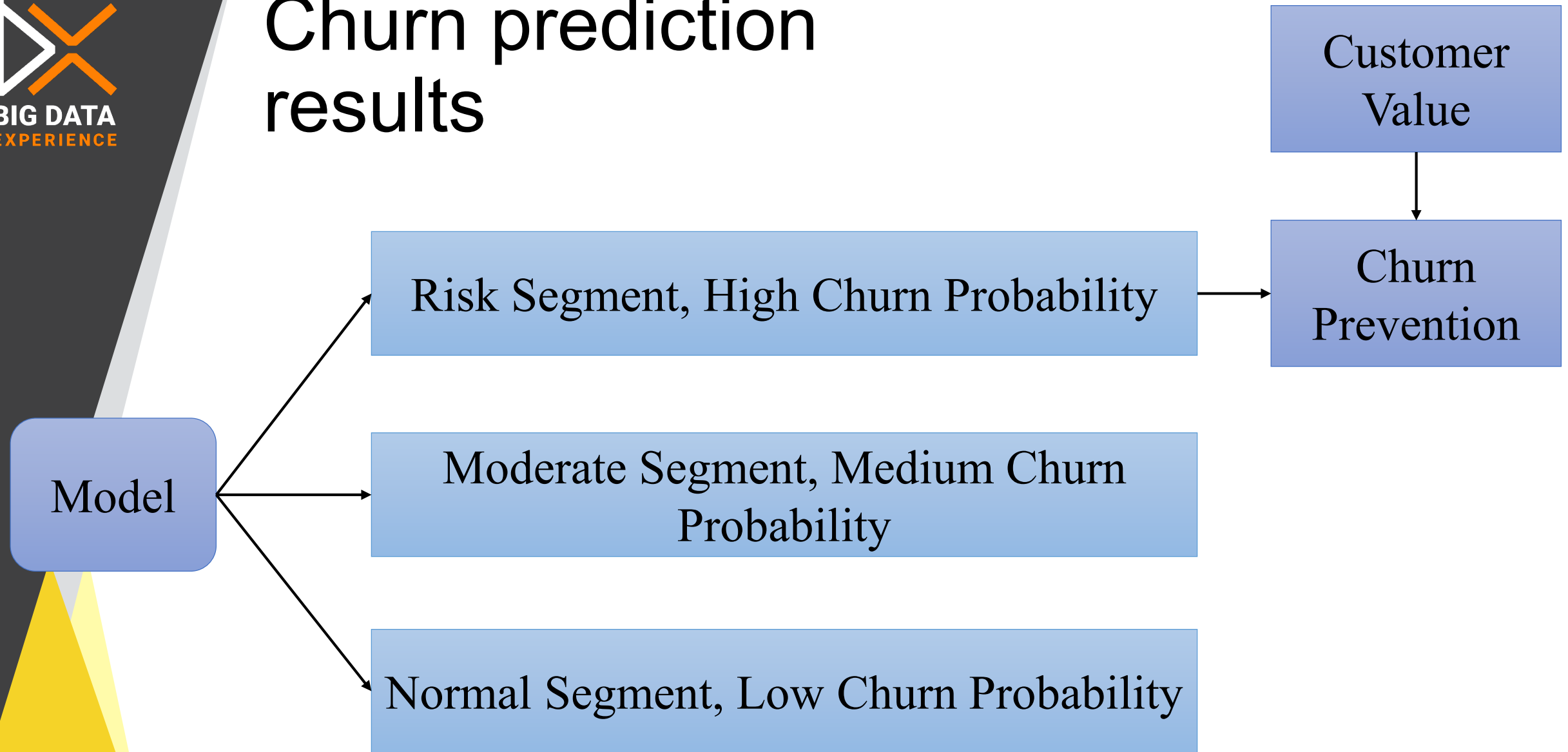
CDR
Payment
Internet usage

Call to contact center
Access to MyAIS

Churn prediction Model



Churn prediction results



Churn prediction Outcome and usage

Outcome

- Able to identify potential churners

Usage

- Offer potential churners with retention campaigns

Data upsell Problem

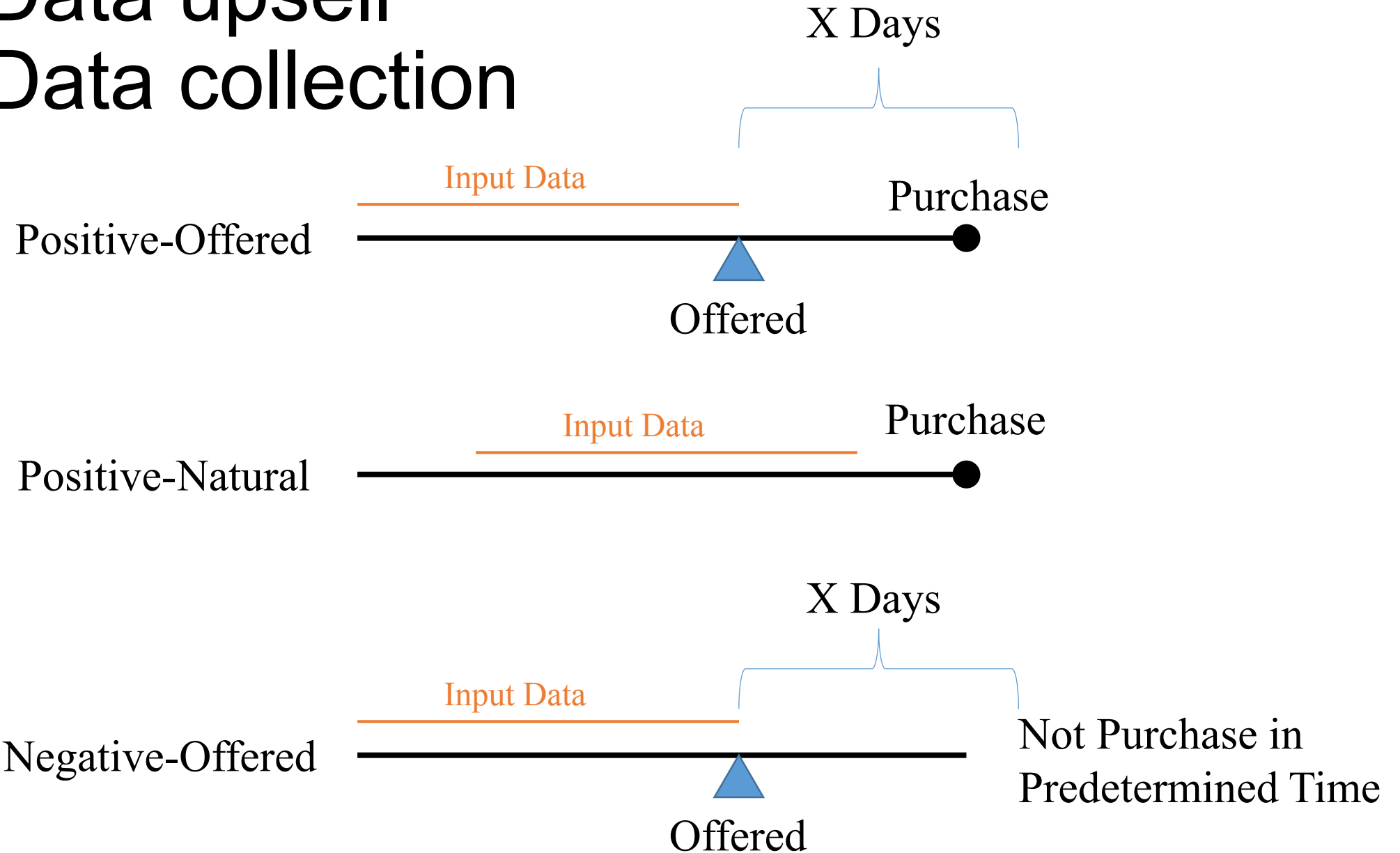


Data upsell Analytic objective

What product to offer? And to whom?

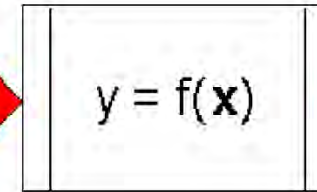
Data upsell

Data collection

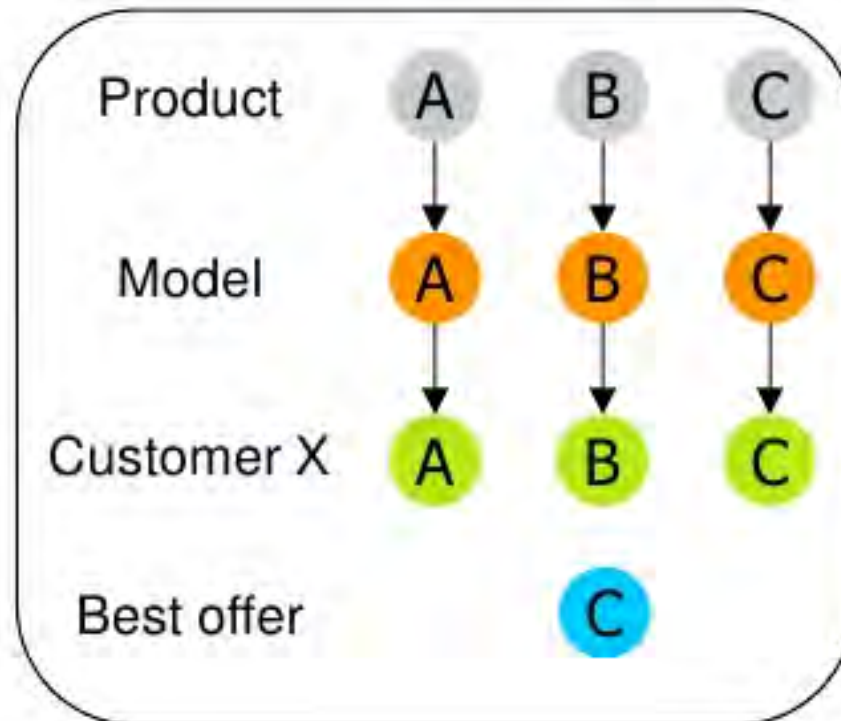


Data upsell Modeling

Predictors	Response
	?
	?
	?
	?
	?
	?
	?



	Propensity Score



Each model is a binary classification model to predict product propensity.

Data upsell

Outcome and usage

Usage

- Connect with the right channel to make automatic offers
- Know which products that each customer are likely to purchase

Outcome

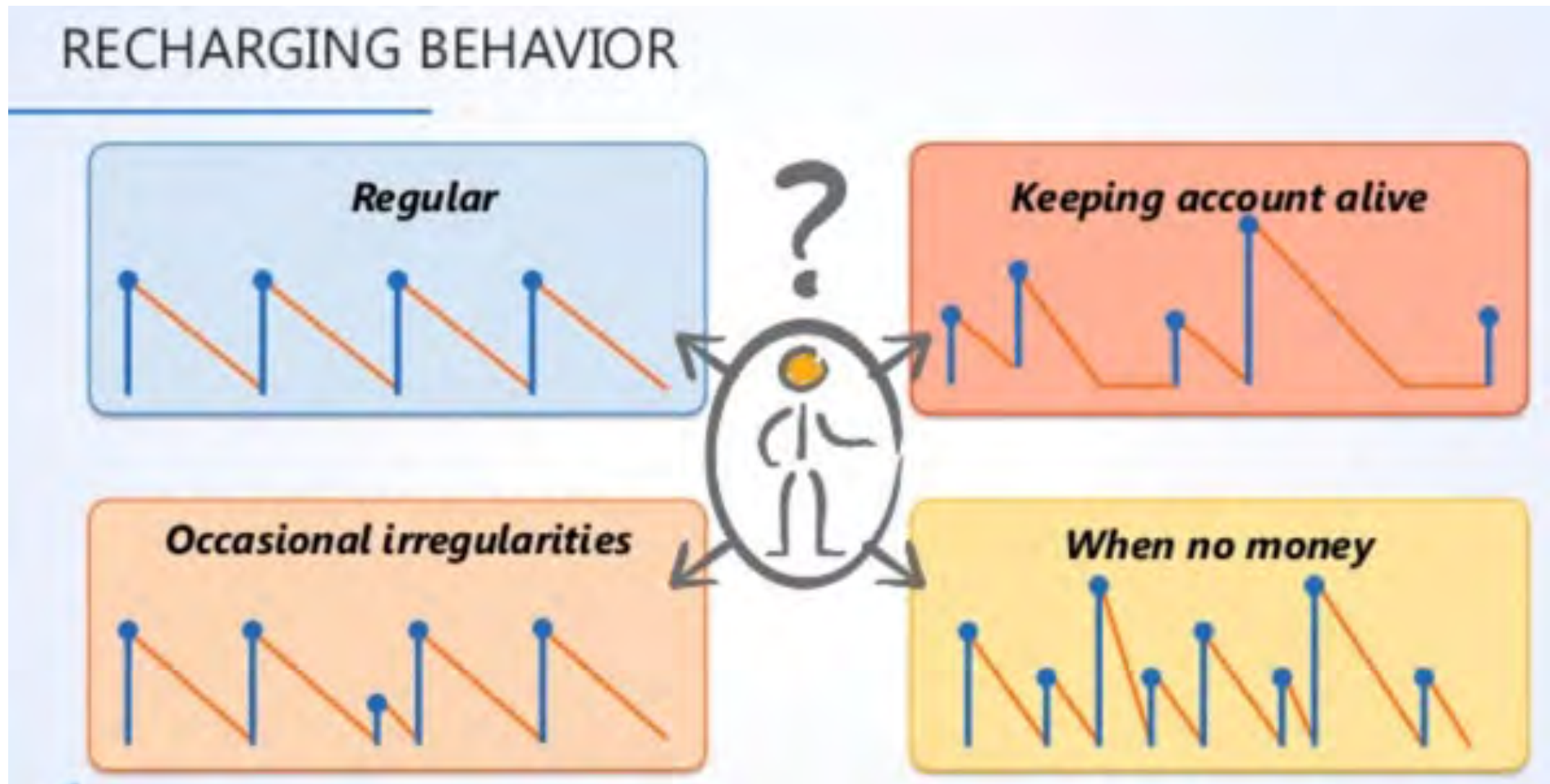
- Increase revenue through automatic upsell

Top-up pattern analysis Problem



Top-up pattern analysis

Top-up behaviors



Top-up pattern analysis Baseline

BEHAVIORAL-DEMOGRAPHIC SEGMENTATION

- survey / usage / demographic data

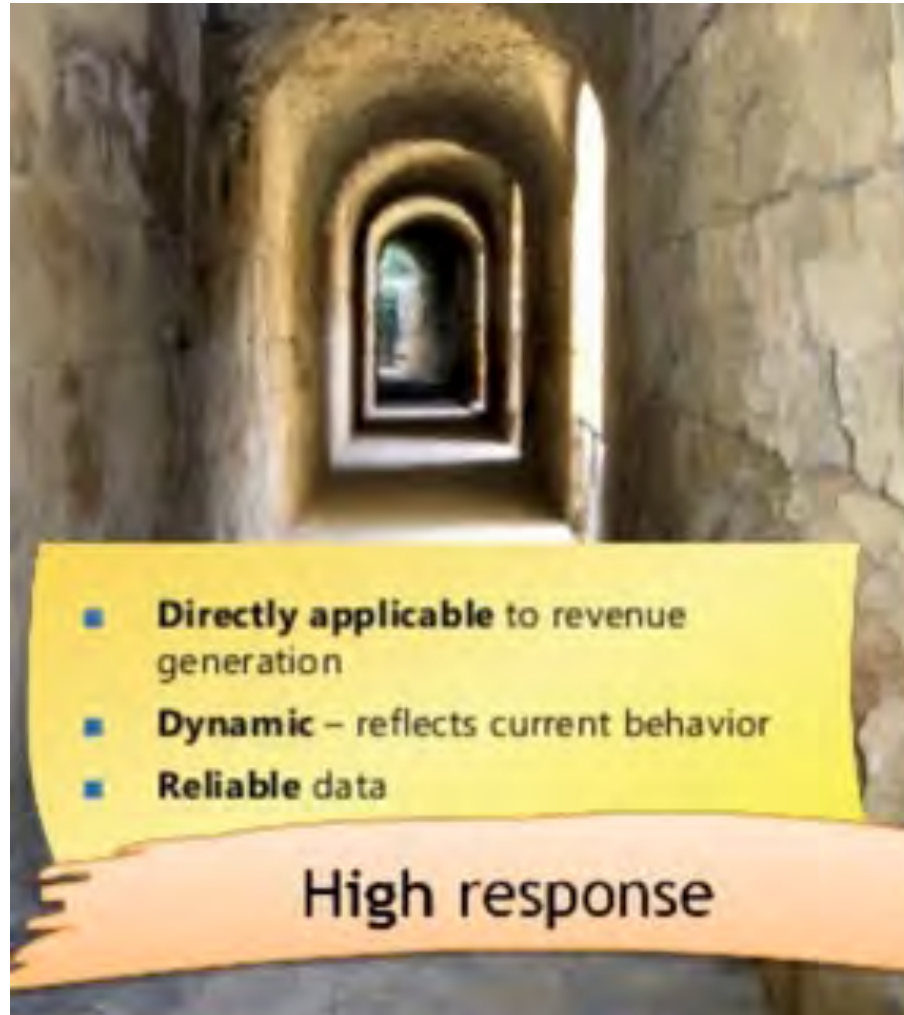


- strategic – for overview
- **hardly applicable** segmentation - mapping surveys to population
- **no direct** link
- **static**
- unreliable & weak data coverage

too general offers, **low response**

Top-up pattern analysis

Expected outcome

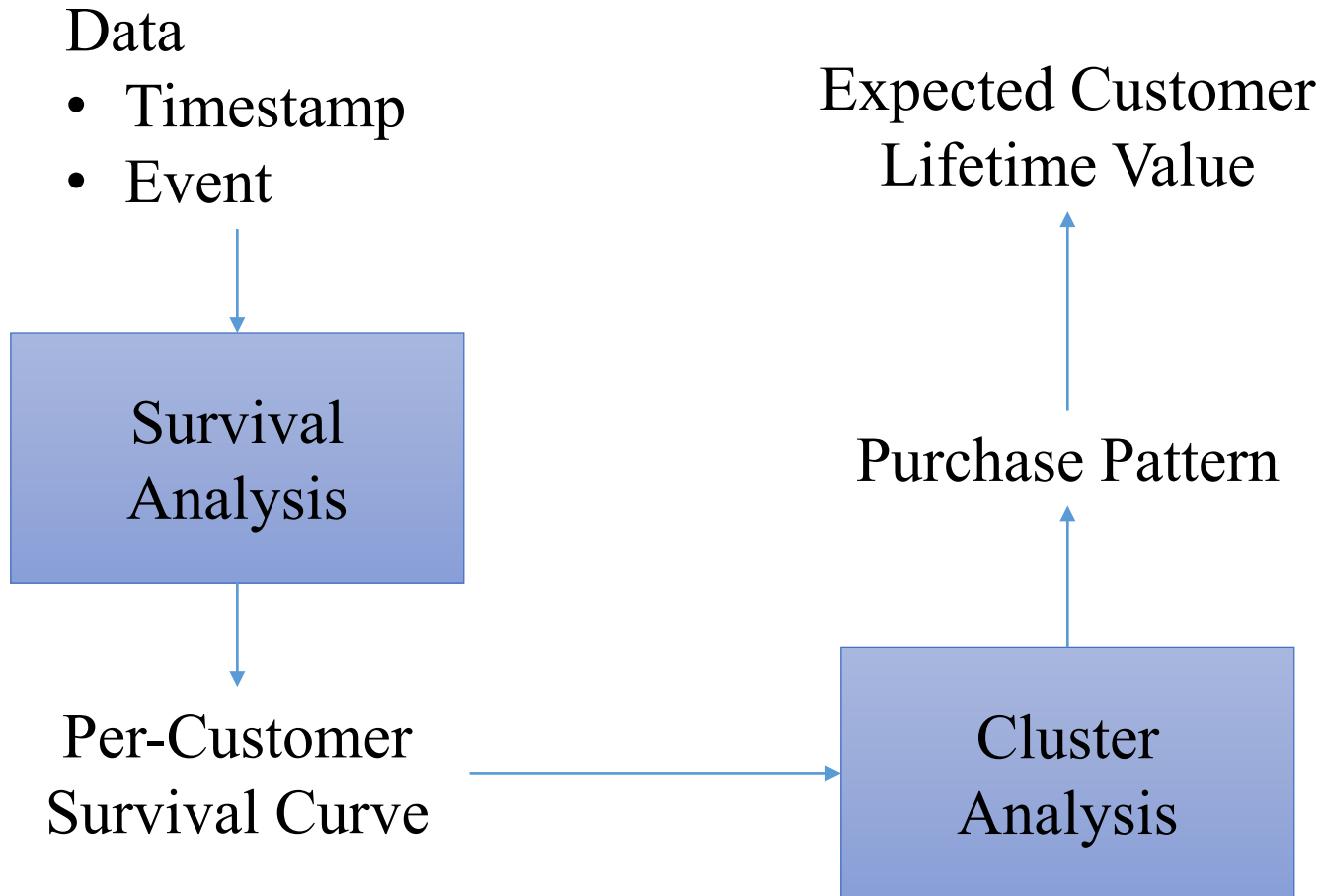


RECHARGE-BASED SEGMENTATION

- based on reliable recharge data
- triggered by customer actions



Modeling: Purchase Pattern



Top-up pattern analysis Marketing plan

PROBLEM & SOLUTION

? How to approach Prepaid users?

do segmentation based
on recharging patterns



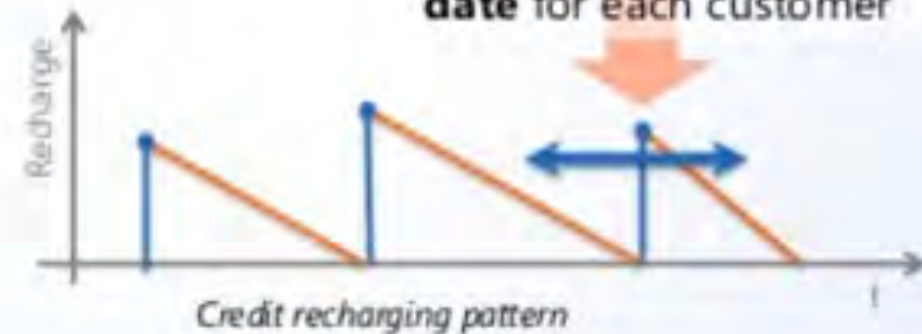
adapt message
to recharge segment



send timely
marketing message

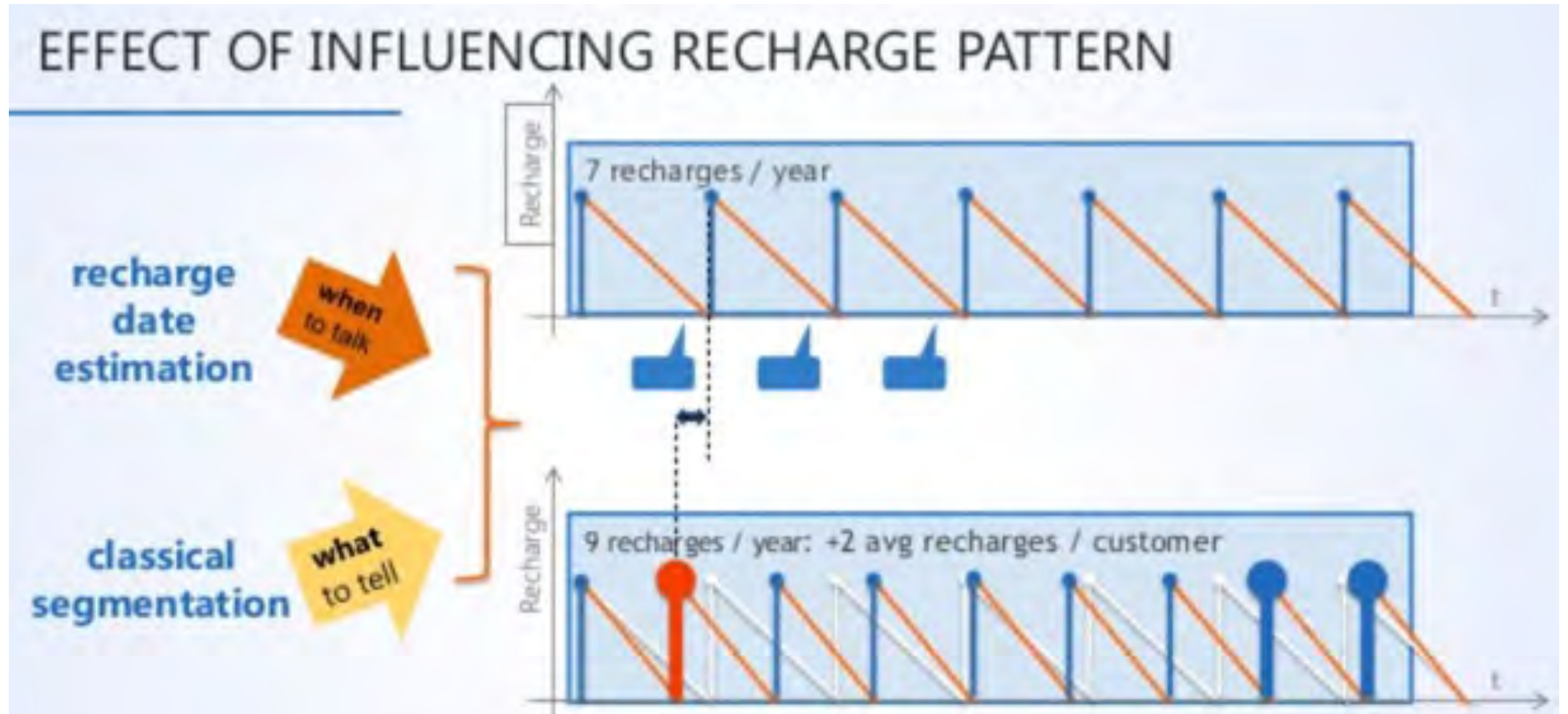
Predictive
models

Estimated recharge
date for each customer



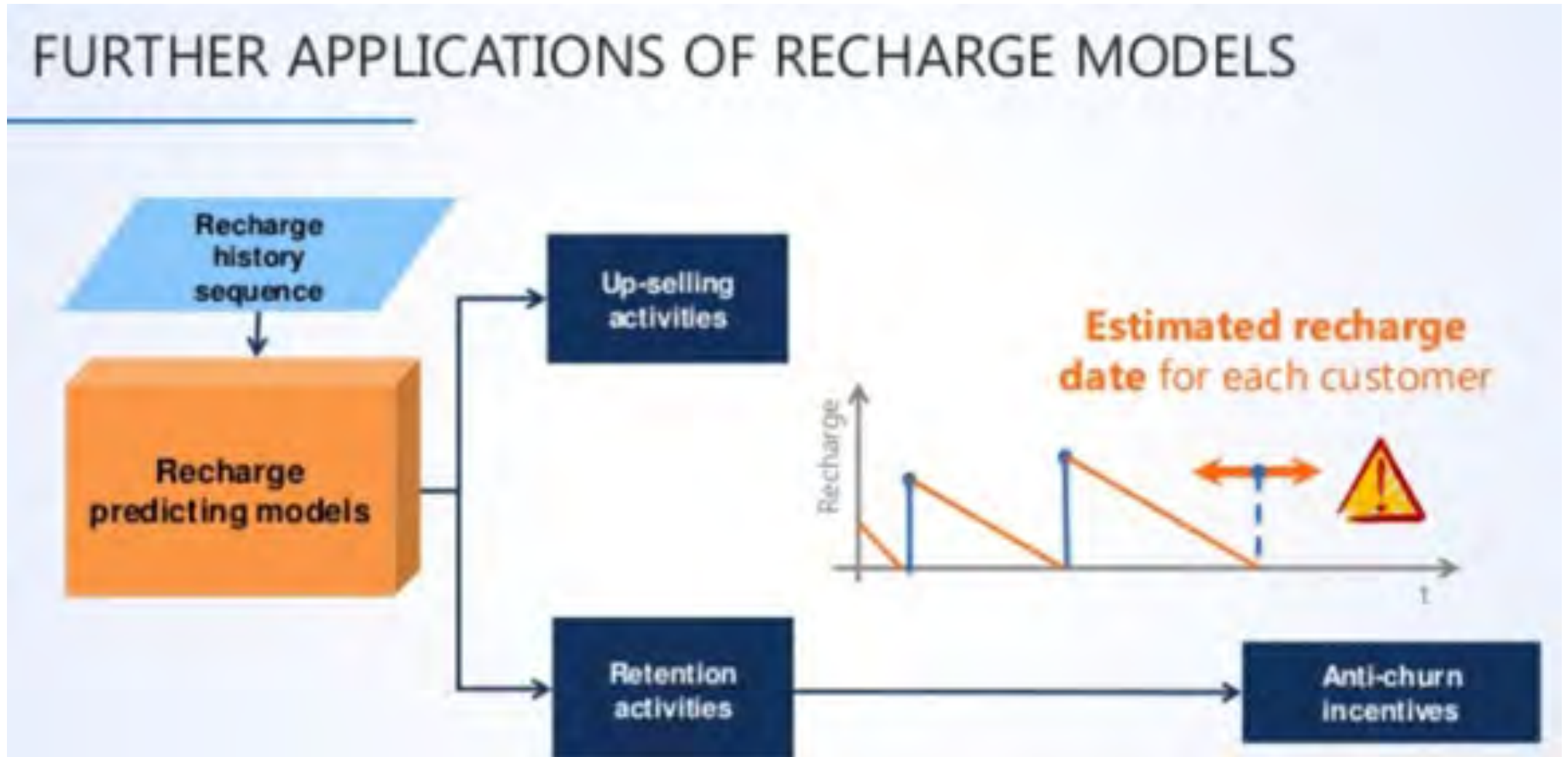
Top-up pattern analysis

Potential influence



Top-up pattern analysis

Further application



Credit risk score



หากคะแนนผ่านเกณฑ์ขั้นต่ำ (Cut-off score) และไม่ขัดกับนโยบายสินเชื่อ (Product policy)



NCB หมายถึง บริษัทข้อมูลเครดิตแห่งชาติ

- ลูกหนี้รายย่อยไม่มีงบการเงินเหมือนลูกหนี้นิติบุคคลจึงอาจไม่มีข้อมูลแหล่งที่มาของรายได้ที่เป็นปัจจุบัน
- จึงใช้ข้อมูล "พฤติกรรม" เพื่อวัดความสามารถในการชำระหนี้

- ติดตามความเสี่ยงลูกค้าแต่ละกลุ่ม
- ใช้คะแนนประกอบการต่ออายุ/วงเงินสินเชื่อ กำหนดอัตราดอกเบี้ย หรืออนุมัติสินเชื่อใหม่ (Product cross-selling)
- คะแนนต่างกัน Action ต่างกัน

Credit risk score Data collection

1.2 การจัดเก็บข้อมูล: เตรียมฐานข้อมูลปัจจัยบ่งชี้ความน่าจะเป็นในการชำระหนี้คืน

ตัวอย่าง

Credit Score

☒ Excellent
☐ Good
☐ Fair
☐ Uncertain

**Ability to pay +
Willingness to pay**

ข้อมูลผู้ขอสินเชื่อ (Demography) มาจากใบคำขอสินเชื่อ

- เพศ อายุ การศึกษา
- อาชีพ / ประสบการณ์ทำงาน
- รายได้ปัจจุบัน

ข้อมูลประวัติการชำระหนี้ (Payment behavior)

- จำนวนครั้งที่ค้างชำระ 12 เดือนล่าสุด
- % การใช้งานเงินเฉลี่ยใน 3 เดือน
- ระยะเวลาไม่ชำระหนี้ใน 6 เดือน
- จำนวนบัตรเครดิตที่เปิดใหม่ใน 6 เดือน
- ยอดหนี้คงค้างทั้งหมด / รายได้
- จำนวนครั้งที่เช็คข้อมูล NCB ในอดีต 12 เดือน

เงื่อนไขการกู้ยืม

- สัดส่วน down payment
- ระยะเวลาการกู้ยืม
- ฯลฯ



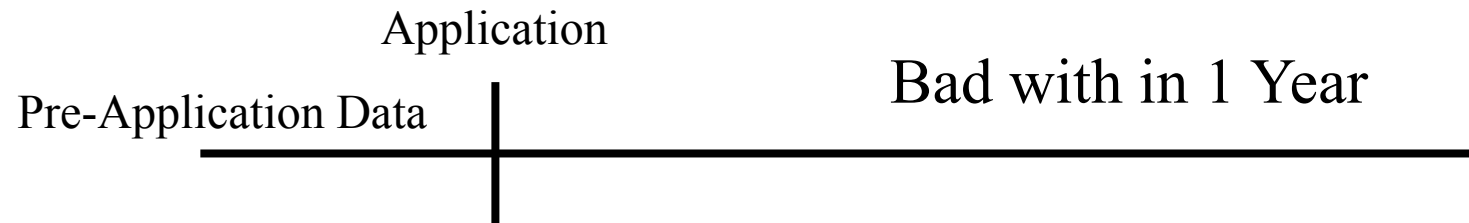
ตัวอย่าง

ID	Gender	Home	...
RB000000000001	F	BKK	...
DG0000000010166	M	Chiang Mai	...

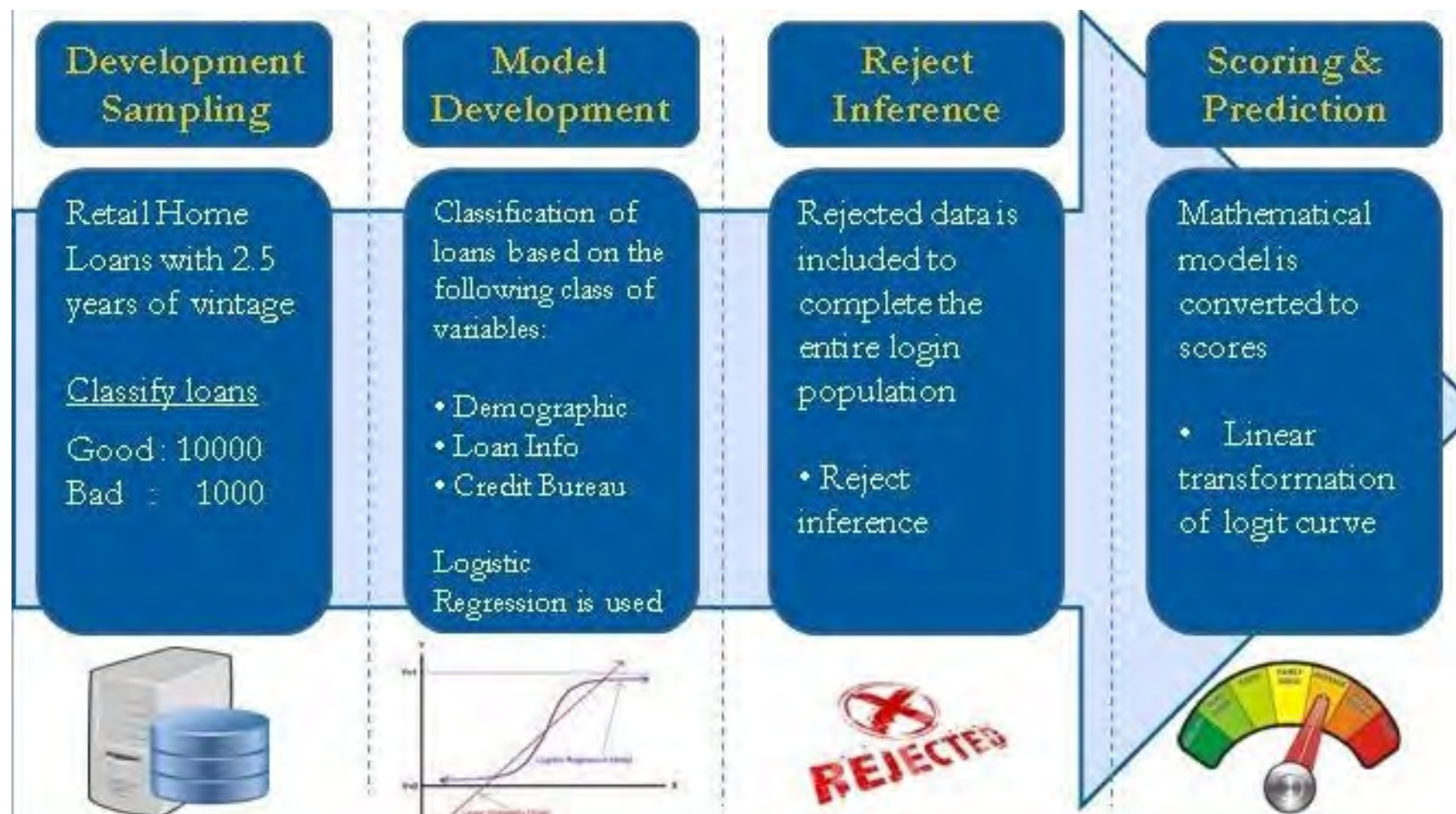
ข้อควรระวัง!

ข้อมูลที่นำมาใช้จัดทำ Credit scoring ต้องไม่สามารถระบุ
ตัวตนของเจ้าของข้อมูลได้

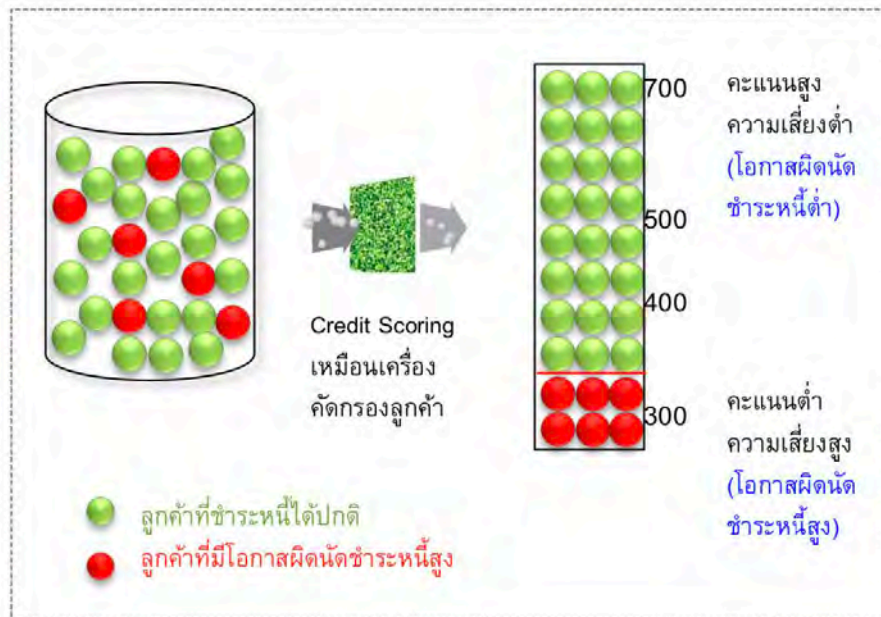
Credit risk score Timeline



Credit risk score Model development



Credit risk score Usage and outcome



ธนาคารพาณิชย์ และสถาบันการเงินต่าง ๆ
จึงใช้ **Credit Scoring** เป็นเครื่องมือประกอบ
การวิเคราะห์สินเชื่อ และอนุมัติสินเชื่อ
โดยเฉพาะสินเชื่อรายย่อย เช่น สินเชื่อ
บัตรเครดิต สินเชื่อบุคคล สินเชื่อบ้าน
สินเชื่อเช่าซื้อรถยนต์ เป็นต้น

Credit scoring ช่วย
เรียงลำดับความเสี่ยงให้
เป็นคะแนนที่เข้าใจง่าย

POOR
300-619

FAIR
620-679

GOOD
680-730

GREAT
730+

Use Case: Product Recommendation



Use Case: Customer Preference

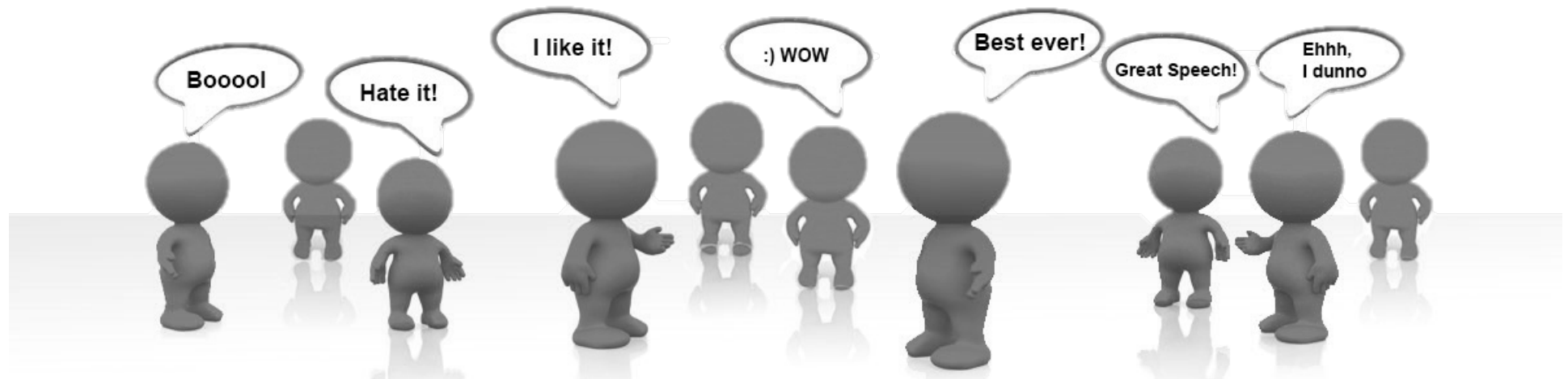
- Zarola derives customer preference and styles based on their transactions
- It optimizes market strategies based on each user profile.



ZALORA

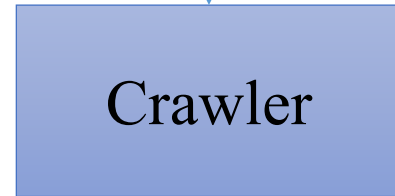
Use Case: Sentiment Analysis

- To determine the attitude of a writer with respect to some topic or the overall contextual polarity of a document.
- Widely applied to reviews and social media for a variety of applications, ranging from marketing to customer service



Use Case: Sentiment Analysis

Social Network Data

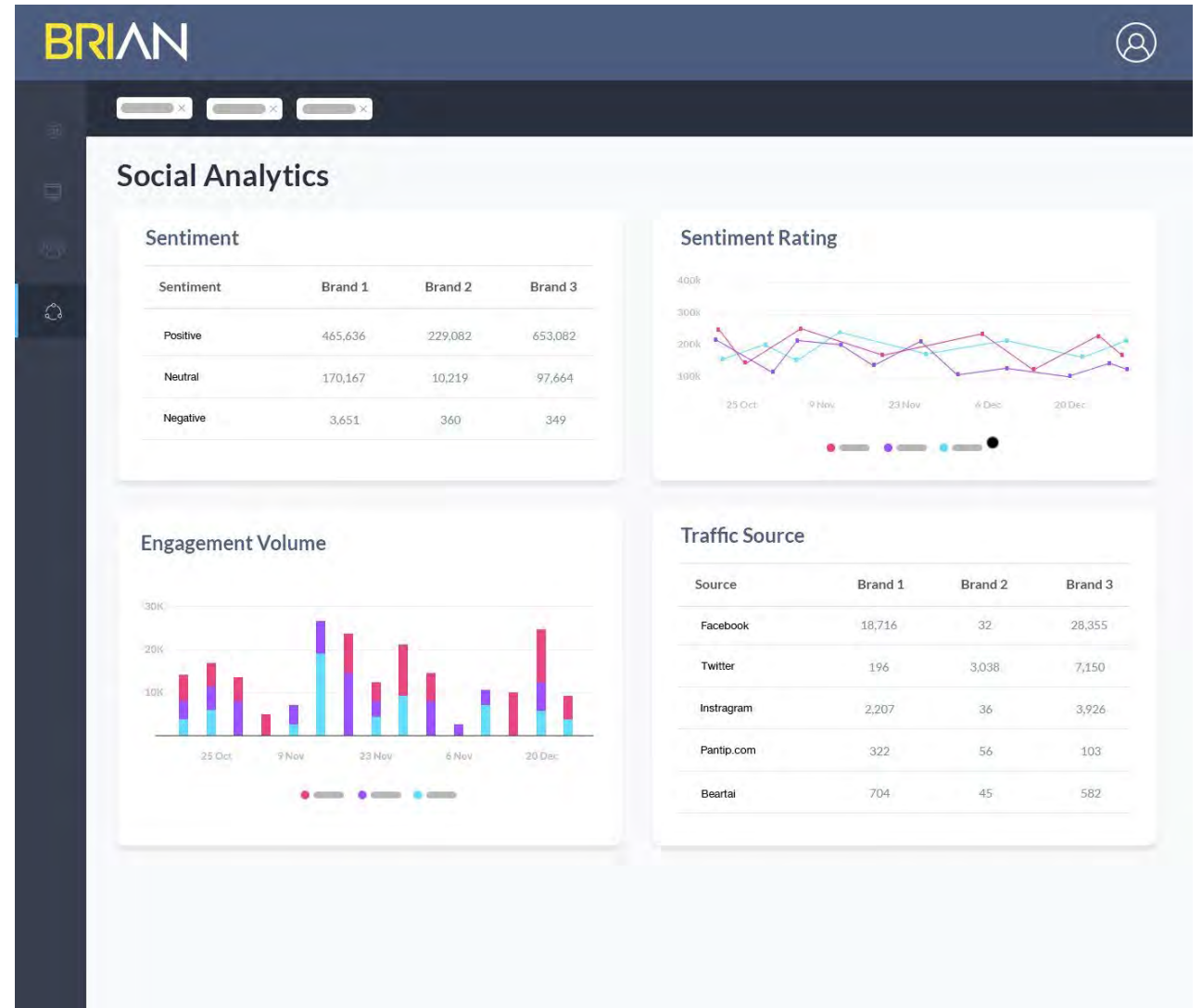
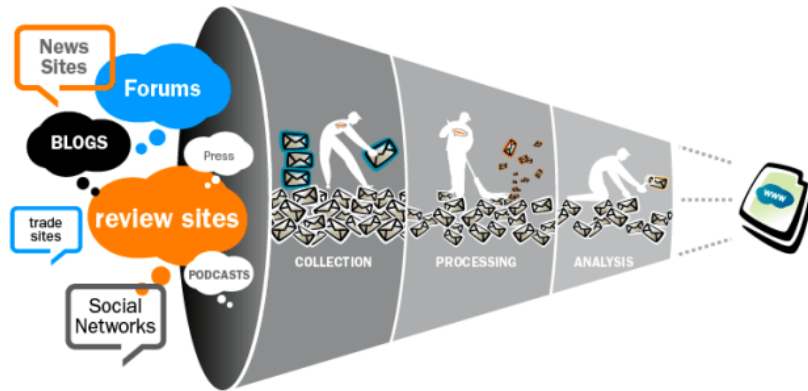


Text Data

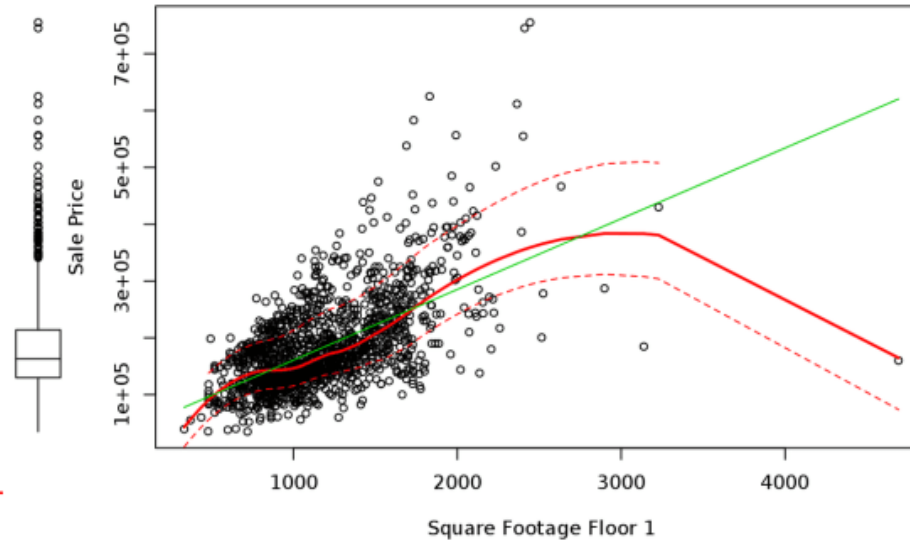
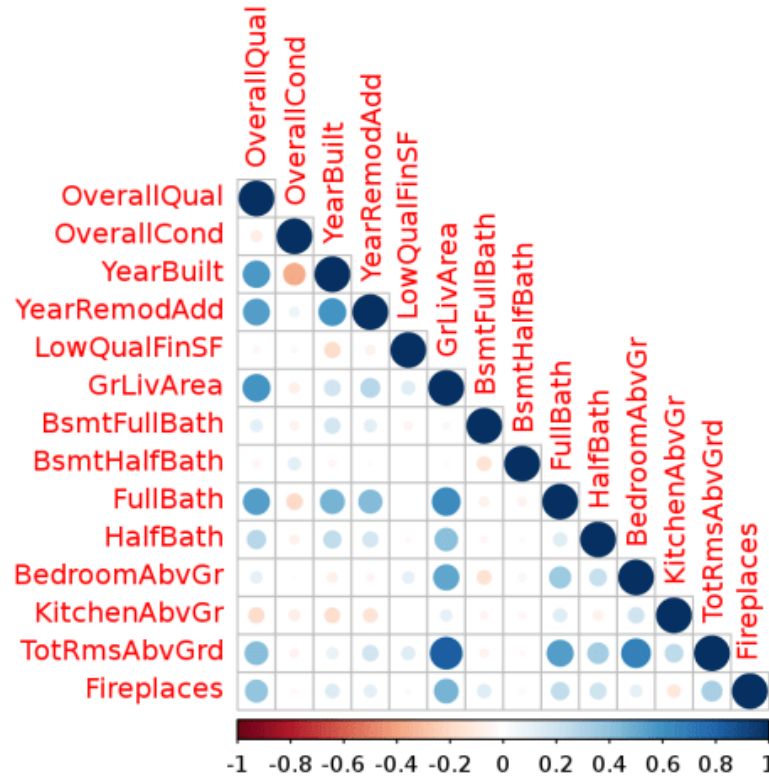
Text + Sentiment



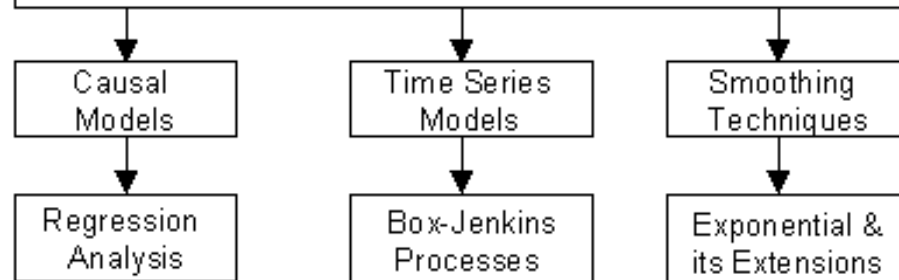
Use Case: Brand Tracking



House Price Prediction



Classification of the Widely Used Forecasting Techniques



End of Lecture 1

Question?