


**NANYANG  
TECHNOLOGICAL  
UNIVERSITY**  
SINGAPORE

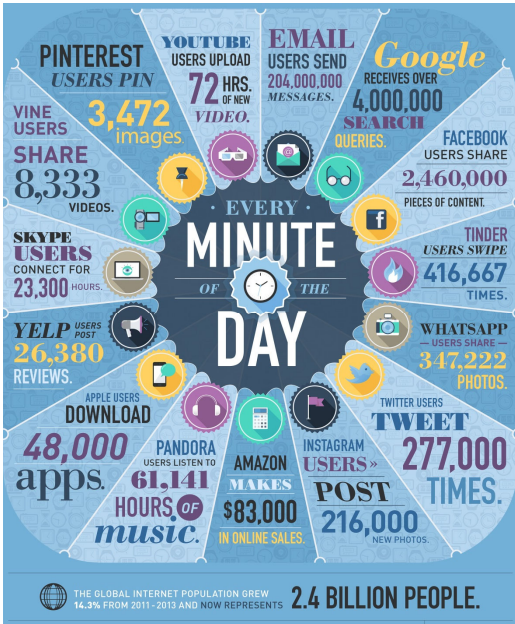
# BC3402

## Weeks 10 to 11

Assoc Prof Goh Kim Huat  
Nanyang Business School



1




**EVERY MINUTE OF THE DAY**

- PINTEREST**: USERS PIN 3,472 Images
- VINE**: USERS SHARE 8,333 VIDEOS.
- SKYPE**: USERS CONNECT FOR 23,300 HOURS.
- YELP**: USERS POST 26,380 REVIEWS.
- APPLE**: USERS DOWNLOAD 48,000 apps.
- PANDORA**: USERS LISTEN TO 61,141 HOURS OF music.
- AMAZON**: MAKES \$83,000 IN ONLINE SALES.
- INSTAGRAM**: USERS POST 216,000 NEW PHOTOS.
- TWITTER**: USERS TWEET 277,000 TIMES.
- WHATSAPP**: USERS SHARE 347,222 PHOTOS.
- TINDER**: USERS SWIPE 416,667 TIMES.
- FACEBOOK**: USERS SHARE 2,460,000 PIECES OF CONTENT.
- Google**: RECEIVES OVER 4,000,000 SEARCH QUERIES.
- EMAIL**: USERS SEND 204,000,000 MESSAGES.
- YOUTUBE**: USERS UPLOAD 72 HRS. OF NEW VIDEO.

THE GLOBAL INTERNET POPULATION GREW 14.3% FROM 2011-2013 AND NOW REPRESENTS **2.4 BILLION PEOPLE.**

Make a guess..... What percentage of the world's data was created in the last 2 years?



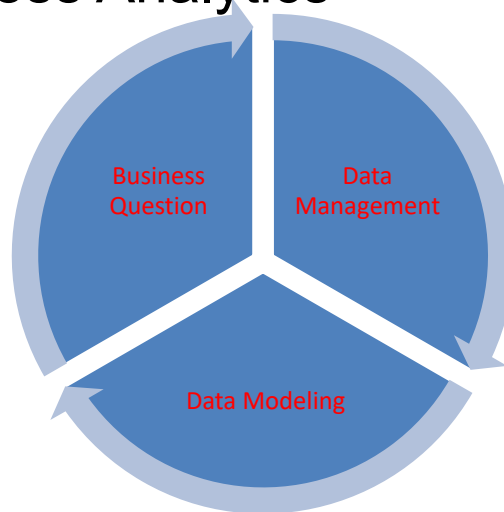
**NANYANG TECHNOLOGICAL UNIVERSITY** SINGAPORE

2



3

## Business Analytics



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

5

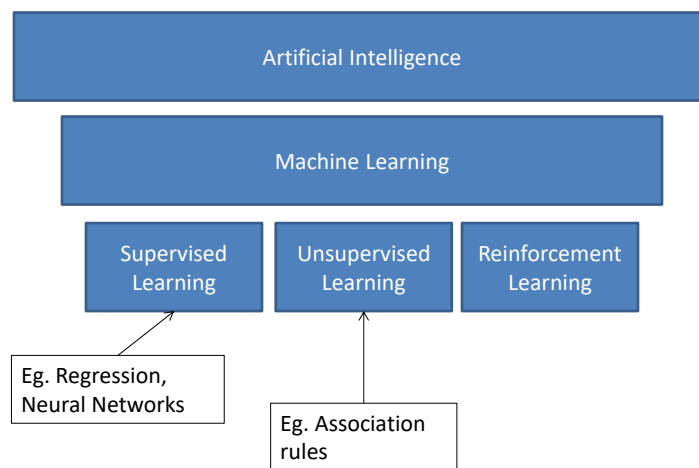
Which of the following is artificial intelligence?



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

6

## Semantic Tree



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

7

## Supervised vs. Unsupervised Alice in Wonderland



"Would you tell me, please, which way I ought to go from here?"  
"That depends a good deal on where you want to get to."  
"I don't much care where –"  
"Then it doesn't matter which way you go."  
- Lewis Carroll



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

8

## Supervised Learning

- Goal: Predict a single "target" or "outcome" variable
- Training data, where target value is known
- Score to data where value is not known
- Uses:
  - Explaining
  - Predicting



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

9

# Unsupervised Learning

- Goal: no specific “target” or “outcome”
- There is no target (outcome) variable to predict or classify
- Broad, fuzzy goals **NOT** specific “targets”
- Multiple purposes:
  - Segment data into meaningful segments
  - detect patterns
  - explore the data
  - etc.....



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

10

# Supervised: Classification

- Goal: Predict categorical target (outcome) variable
- Examples: Purchase/no purchase, fraud/no fraud, creditworthy/not creditworthy...
- Each row is a case (customer, tax return, applicant)
- Each column is a variable
- Target variable is often binary (yes/no)



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

11

## Supervised: Prediction

- Goal: Predict numerical target (outcome) variable
- Examples: sales, revenue, performance
- As in classification:
- Each row is a case (customer, tax return, applicant)
- Each column is a variable
- Taken together, classification and prediction constitute “predictive analytics”



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

12

## Unsupervised: Association Rules

- Goal: Produce rules that define “what goes with what”
- Example: “If X was purchased, Y was also purchased”
- Rows are transactions
- Used in recommender systems – “Our records show you bought X, you may also like Y”
- Also called “affinity analysis”



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

13

## Unsupervised: Data Reduction

- Distillation of complex/large data into simpler/smaller data
- Reducing the number of variables/columns (e.g., principal components)
- Reducing the number of records/rows (e.g., clustering)



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

14

## Unsupervised: Data Visualization

- Graphs and plots of data
- Histograms, boxplots, bar charts, scatterplots
- Especially useful to examine relationships between pairs of variables



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

15

## Data Exploration

- Data sets are typically large, complex & messy
- Need to review the data to help refine the task
- Use techniques of Reduction and Visualization



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

16

## Before we start: Data Scoping

- Unit of analysis problems
- Datafication
- Prediction of success problems

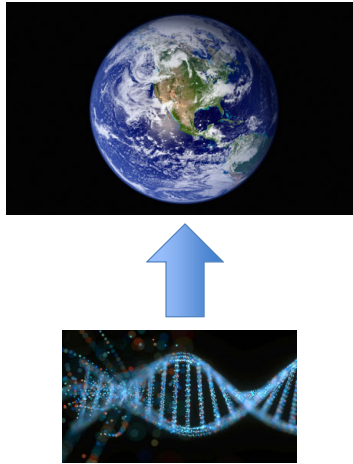


NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

17



## Unit of Analysis



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

18

## Datafication: Leveraging on Existing Data



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

27

## Matthew Fontaine Maury: Pathfinder of the sea

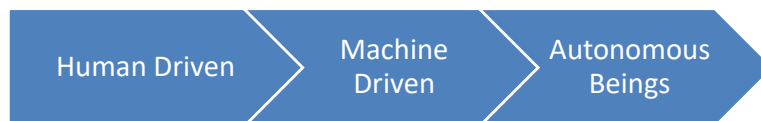


NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

28

## Three Phases

- Consumer
- Business
- Image recognition
- Video analytics
- Text analytics
- Drones
- Cars
- Machinery



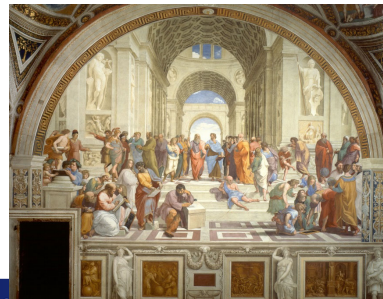
**Datafication**



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

30

## Renaissance [1400 – 1600]



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

31

## (Post) Impressionism [1865 – post 1900s]



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

32

# Concept of Error: Survey

LTA survey findings on MRT commuter graciousness greeted with scepticism

Published on  
Aug 22, 2013  
6:50 AM

(http://www.asias  
www.asias  
straitsnews.com)  
pres.co  
f  
t  
s



Will you give up your seat to someone who needs it more? **94% Yes**

Will you queue up and give way to arriving passenger? **98% Yes**

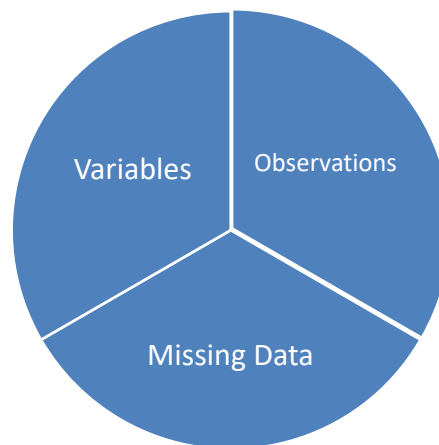
Will you move in for others to board the train? **96% Yes**



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

33

# Predicting Success of Projects



Maccabitech  
• 3V, 4V, 5V?  
• NDA?



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

34

SUPERVISED LEARNING

## TECHNIQUES TO EXPLAIN AND TO PREDICT



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

35

## The Hierarchy

Level 4 – Prescriptive Analytics

Suggest what you should be doing after weighing the possible states

Level 3 – Predictive Analytics

What could happen? Different possible states

Level 2 – Diagnostic Analytics

Explaining what is happening

Level 1 – Descriptive analytics

What happened



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

36

## Uses of Supervised Learning Tools

### Explanation

- Give reasons why?
- Policy-making/  
decision-making/  
Justification
- Business/ consumer  
insights
- “More” scientific

### Prediction

- If X, then Y
- More application in  
operations
- May not care about  
the underlying cause
- May be “less”  
scientific



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

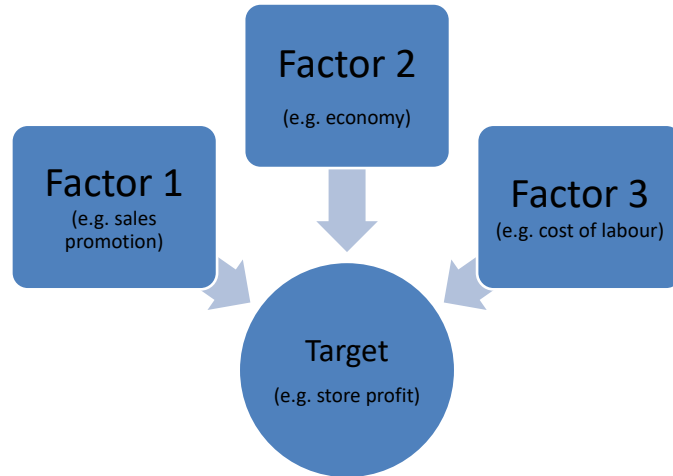
37



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

38

## Regression (To Explain)



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

39

Multiple Regression: **EXPLAINING**

## Traffic Statistics (Singapore)

	2010	2011	2012
Speeding Violations	205,623	225,550	244,806
Red-Running	17,185	17,492	17,705
All traffic Violations	304,472	316,214	327,503

Are Singapore drivers becoming more reckless?

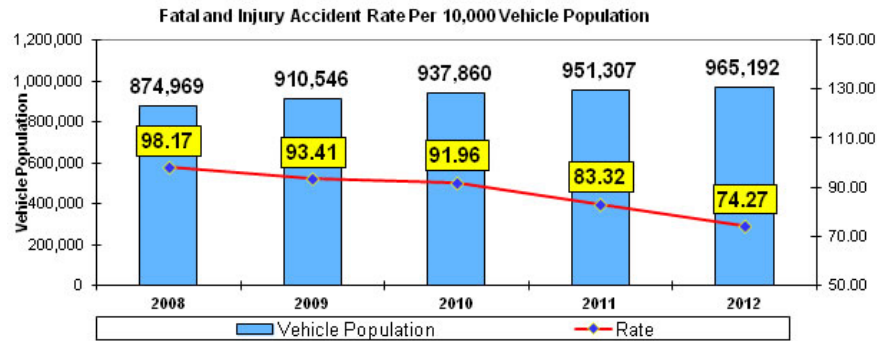
What is the “target” for this model to answer the above question?



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE  
<http://www.ntu.edu.sg/home/linsheng/>

40

## Traffic Statistics (Singapore)



Are Singapore drivers becoming more reckless?

What is the "target" for this model to answer the above question?



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

41

## Proxies & Partial Out The Reasons

- Target measure is never perfect. Proxies used
  - Accuracy of proxy
- Need other factors to partial out the alternative explanations
- More factors greater fidelity

*"Holding XYZ constant"  
here are my conclusions*



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

42



Questions to (Transaction) Data

## CASES OF SUPERVISED LEARNING



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

45

## ATM

### Banks get a peek into customers' whims and needs Big Data analytics provides unprecedented insights, prompts change

BY AMIT ROY CHOUDHURY

[PRINT](#) | [EMAIL THIS ARTICLE](#)



Real-time reach: Banks can engage customers based on where they are and what they are doing. - FILE PHOTO

[SINGAPORE] How often should cash at ATMs be replenished? Which products should specific customers be offered?

Thanks to Big Data analytics, Singapore-based banks are taking the guesswork out of such questions, while gaining deeper insights into customer preferences.

David Gledhill, managing director and head of group technology and operations at DBS, said that the bank used Big Data analytics to improve cash availability for its islandwide ATM service while saving on costs.

**Logic + Creativity + Statistics**



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

46

## Three key questions

- Is there “something” that we need to predict or explain to achieve our overall objective
- What is the target for the model
- Can you see the target in the dataset



NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

48

## A Different Type of Transaction...

MARINA BAY   
SINGAPORE

Panopticon



NANYANG TECH

51

# Detecting Fraud



$d$	$P(d)$
1	30.1%
2	17.6%
3	12.5%
4	9.7%
5	7.9%
6	6.7%
7	5.8%
8	5.1%
9	4.6%

$$P(d) = \log_{10} (1+1/d)$$

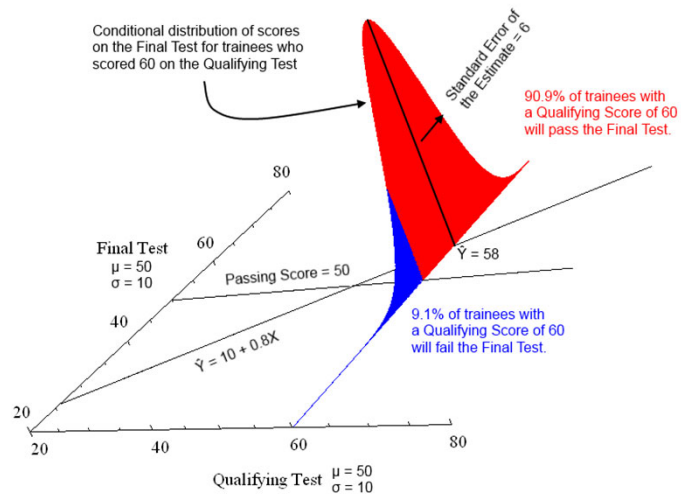
*Genius you are vs. the genius you have: Elizabeth Gilbert*



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

Logic + Creativity + Statistics + Luck

52



It is about predicting a distribution ...



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

53

## Supervised Learning: Predicting Transaction Distribution

- Distribution of multi-dimensional behavior (ATM)
- Distribution of known probability functions (Casino)
- Distribution of natural numbers – outliers detection (Fraud detection)

