

Introduction to Data Science

Dr. Sandareka Wickramanayake

sandarekaw@cse.mrt.ac.lk

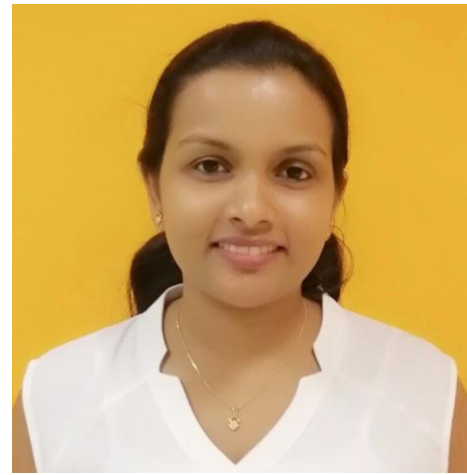
Learning Outcomes

After completing this module, students should be able to

- Demonstrate data acquisition, data representation, and data pre-processing skills to describe, analyze and repurpose data from a variety of sources.
- Apply critical thinking and statistical techniques to understand and visualize relationships in data
- Apply machine-learning techniques in exploratory data analysis for problems related to commerce, industry, and research.
- Design and compute a statistical relationship in data including correlation and linear regression
- Design and develop data-driven algorithms for outcome prediction

Delivery

- Lectures - Mon 10.15 AM – 12.15 PM (CSE Level 2 Lab)
- Labs – Wed 3.15 – 5.15 PM (Sem 4), 3.15 – 5.15 PM (Sem 8)
- Lecturers
 - Dr. Nisansa De Silva - NisansaDdS@cse.mrt.ac.lk
 - Dr. Sandareka Wickramanayake - sandarekaw@cse.mrt.ac.lk



Course Outline

Week	Lecture Topic	Lecturer
1	Introduction	SW
2	Data collecting, data documenting, data quality	SW
3	Data preprocessing	SW
4	Descriptive analysis	SW
5	Exploratory analysis	SW
6	Hypothesis Testing	NdeS
7	Visualization and Dashboarding	SW
8	Project Week	
9	Supervised Learning	NdeS
10	Unsupervised Learning and Evaluation	NdeS
11	Prescriptive and Cognitive Analytics	NdeS
12	Big Data	NdeS
13	Ethics	NdeS
14	Data Science Project Evaluation and Discussion	

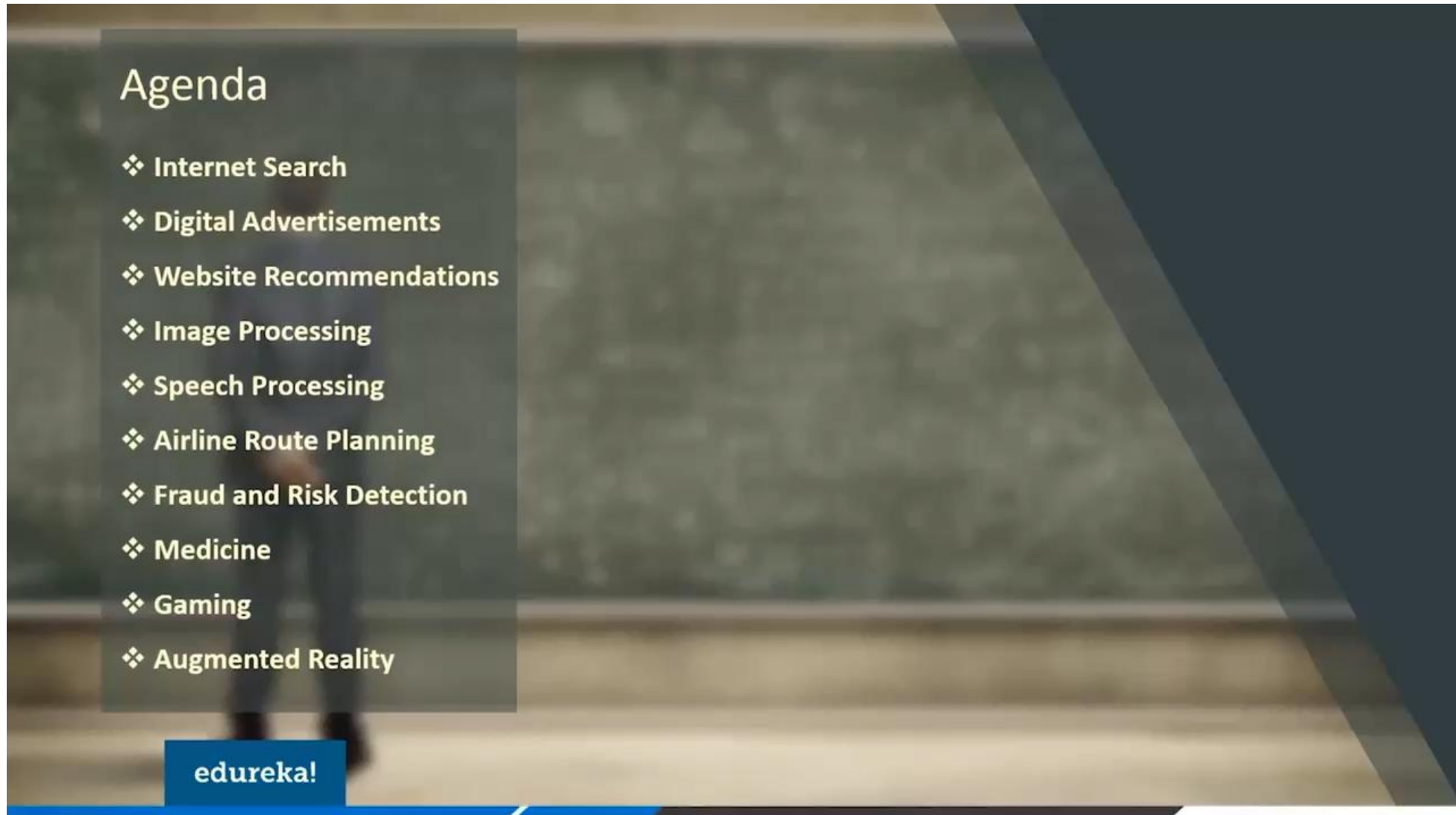
Assessments

- Continuous Assessment – 40%
 - Quizzes (Individual) and in-class activities (Group) – 10%
 - Class project (Group) – 30%
- Final Examination – 60%
 - Online examination conducted in CSE labs
 - 2 hours
 - Open book?

Reading Materials

- No specific textbook
- Additional reading materials related to each topic will be posted on Moodle.

Are We Using Data Science Products?



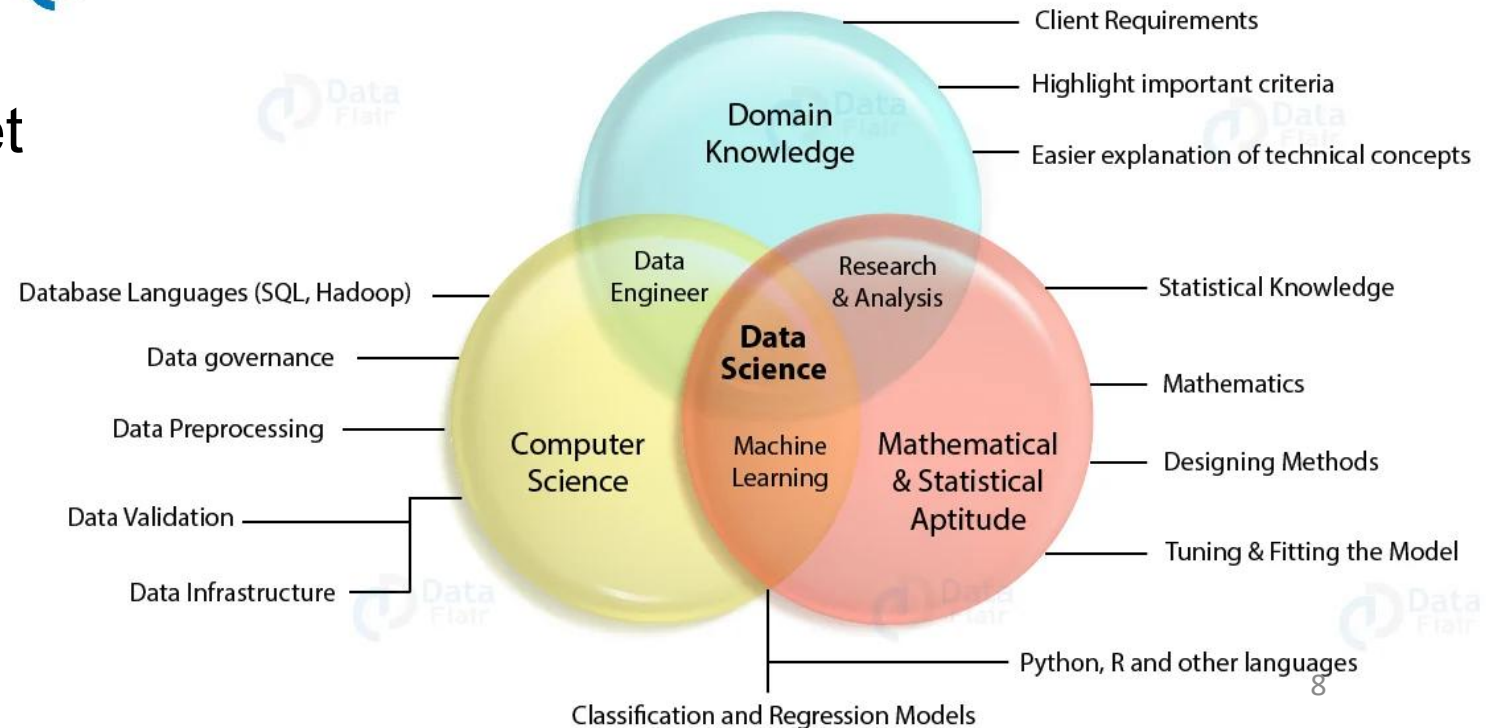
What is Data Science?

- Data Science

- Is the extraction of knowledge from large volumes of data.
- Uncovers actionable insights hidden in data that can be used to guide decision-making and strategic planning.
- Combines many fields.

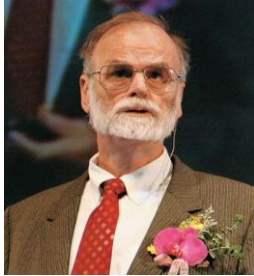


- However, there is not yet a definition agreed upon by all.



What is Data Science?

- Data science = the Fourth Paradigm of Science. – Jim Gray



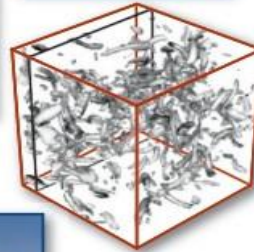
(1942-2012)

Science Paradigms

- Thousand years ago:
science was **empirical**
describing natural phenomena
- Last few hundred years:
theoretical branch
using models, generalizations
- Last few decades:
a **computational** branch
simulating complex phenomena
- Today: **data exploration** (eScience)
unify theory, experiment, and simulation
 - Data captured by instruments or generated by simulator
 - Processed by software
 - Information/knowledge stored in computer
 - Scientist analyzes database/files using data management and statistics



$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G \rho}{3} - K \frac{c^2}{a^2}$$



What is Data Science?

- Data science = Science of data
 - The intellectual and practical activity encompassing the systematic study of facts and statistics collected for reference or analysis.

Google's
definition of
"data"

data

/ˈdɛɪtə/ 

noun

facts and statistics collected together for reference or analysis.

"there is very little data available"

synonyms: facts, figures, **statistics**, details, particulars, specifics, features;

Google's
definition of
"science"

science

/ˈsaɪəns/ 

noun

the intellectual and practical activity encompassing the systematic study of the structure and behaviour of the physical and natural world through observation and experiment.

What is Data Science?

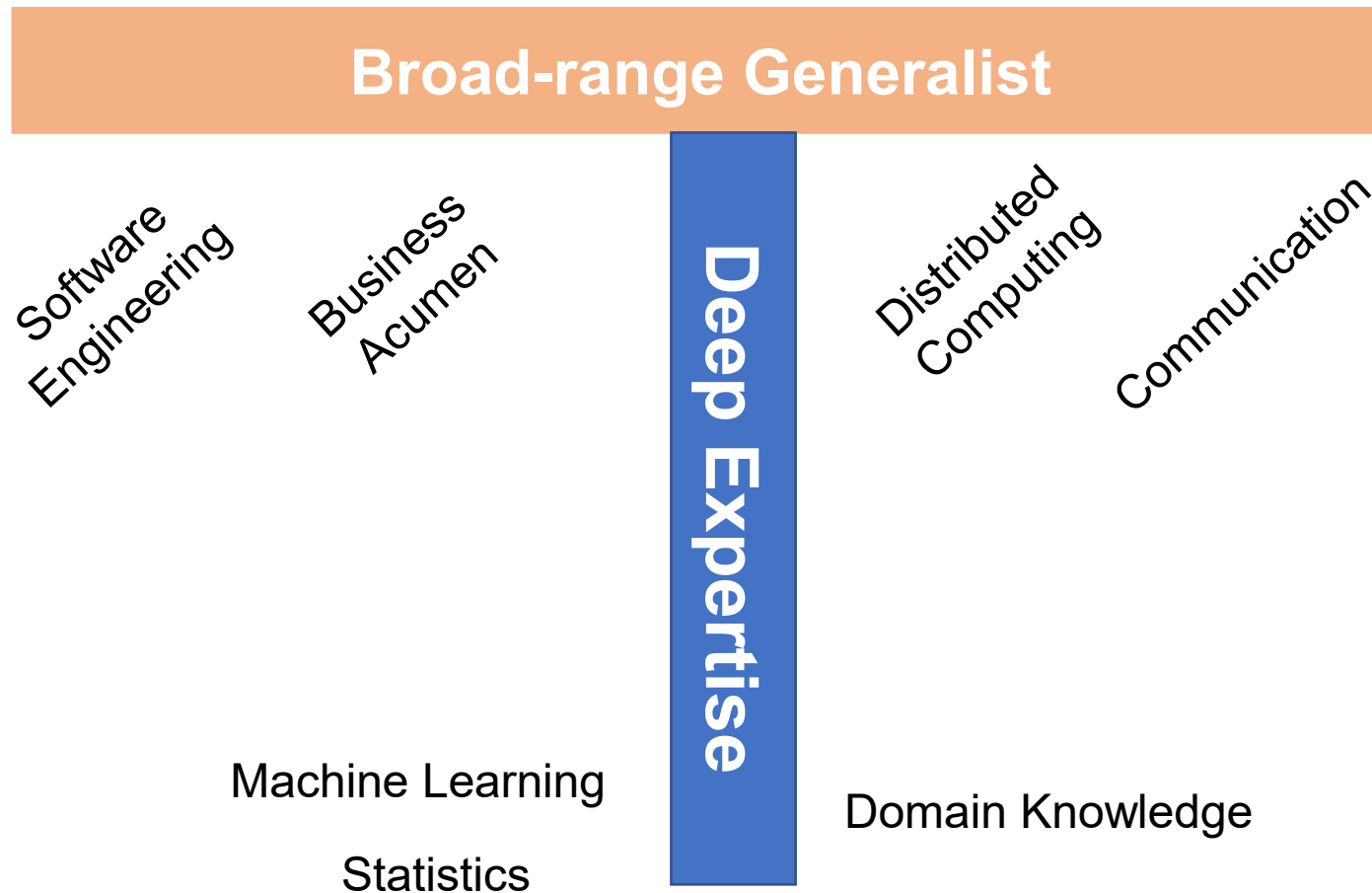
Wikipedia	“Data Science is the extraction of knowledge from large volumes of data that are structured or unstructured”
NIST, 2015	“Data science is the empirical synthesis of actionable knowledge from raw data through the data lifecycle process”
Dhar, 2013	“Data science is the study of generalizable knowledge from data”
Peter Naur, 1974	“[data science is] The science of dealing with data, once they have been established, while the relation of the data to what they represent is delegated to other fields and sciences.”

What is Data Science?

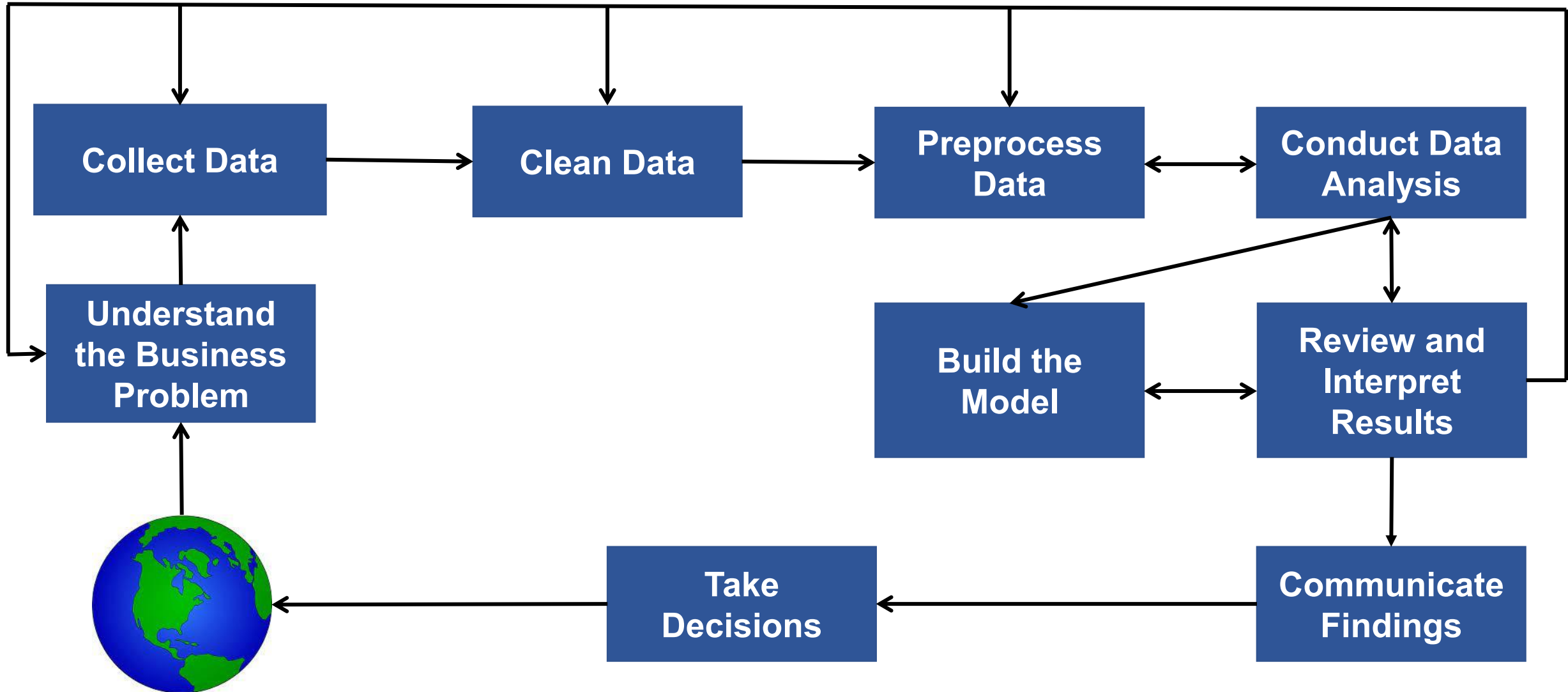
- Data science is an emerging discipline.
 - It remains a science where new knowledge and tools are still being invented.
- There is not yet a clear definition agreed upon by all for the term 'data science'.
 - Different definitions exist from different perspectives (government, business, research, etc.)
 - We adapt NIST's definition: "Data science is the empirical synthesis of actionable knowledge from raw data through the data lifecycle process"
- You, as the future data scientist, will shape the field.

To Succeed in Data Science

- You need the skills of a good Software Engineer and skills in Machine Learning.



Data Science Process



Data Science Process

simplilearn



<https://youtu.be/X3paOmcrTjQ?t=10>

Data Science Process Using a Real World Example



<https://www.youtube.com/watch?v=KdgQvgE3ji4&t=59s>

Data Science Process

- In-class activity 1 - Group activity
 - Check Moodle

Data-Driven Solutions for a Better Sri Lanka

- In-class activity 2 - Group activity
 - The following are 7 recent news headlines.
 1. Tourism earnings target missed by \$1.8 billion in 2025
 2. Lack of coordination between highland institutions makes recovery difficult
 3. Rebuilding Sri Lanka Fund crosses Rs. 8.5bn
 4. Education reforms hit bumps in the road
 5. Sri Lanka's deadly roads claim 1,800 lives in 2025, threatening the tourism industry
 6. Commuters keep complaining: Sri Lanka's public transport sector is trapped in transit
 7. The Exodus and Its Toll: Sri Lanka's Economic Crisis and the Migration of Doctors
 - The topic assigned to your group = your group number % 7 + 1

Data-Driven Solutions for a Better Sri Lanka

- In-class activity 2 - Group activity
 - Your task
 - Brainstorm how DS can be used to address the issue/achieve the target.
 - Identify 2 key applications of DS for the assigned topic.
 - Identify the data sources.
 - What are the anticipated challenges?
 - What is your proposed plan? Discuss how you apply the data science process for the identified applications.
 - Create a two pages report.