

# Collaborative Data Science and COVID-19

Assoc Prof Peter Julian Cayton, PhD

School of Statistics, University of the Philippines Diliman

Member, UP COVID-19 Pandemic Response Team

Member, LEADS for Health Security and Resilience



1

## Flow of the Presentation

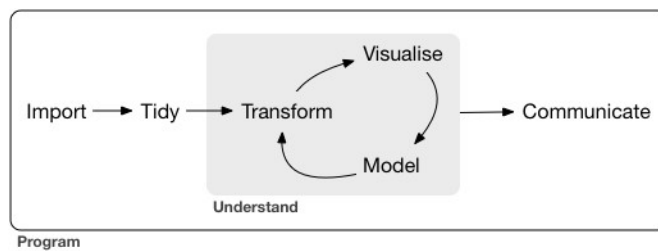
- ▶ A Typical Data Science/Analytics Workflow
- ▶ Membership to Data-Driven Teams
- ▶ Example Projects
- ▶ Demonstrating the Workflow with a Team using Example Projects
- ▶ Closing Remarks

2

## A Typical Data Science/Analytics Workflow

3

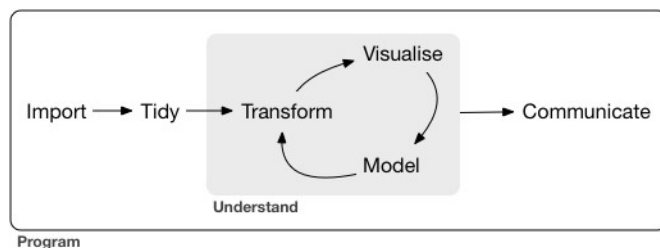
## A Typical Data Science/Analytics Workflow



- Source: Wickham, H and Golemund, G (2017). R for Data Science. O'Reilly. <https://r4ds.had.co.nz/introduction.html>

4

## A Typical Data Science/Analytics Workflow

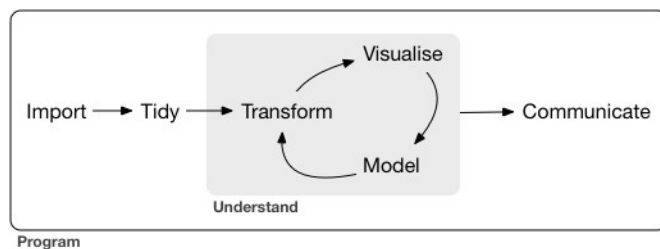


### ► Import

- Extracting the data from an internal database, a file, an online website, or thru a web application programming interface (API), to be loaded in R/RStudio

5

## A Typical Data Science/Analytics Workflow

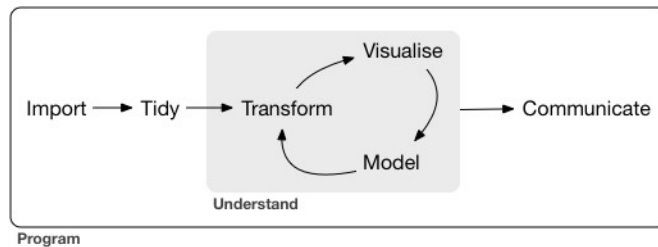


### ► Tidy

- Arranging the data into a neat data structure, with variable as columns and data points as rows.
- Included in this step would be data cleaning, data augmentation, missing data imputation, and others

6

## A Typical Data Science/Analytics Workflow

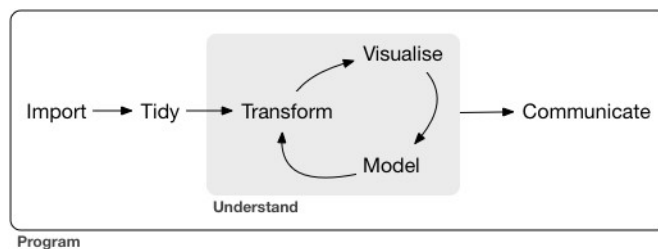


### ► Understand

- Generally, the steps to extract insights from data after tidying up.

7

## A Typical Data Science/Analytics Workflow

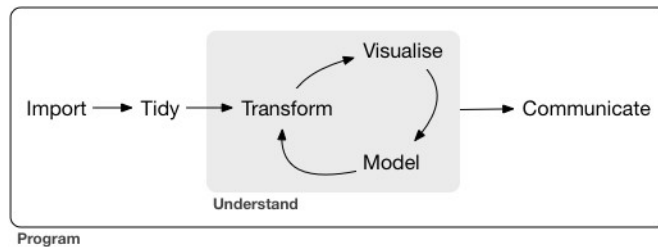


### ► Transform

- Processing the data in preparation for further steps. Examples are:
  - Narrowing the data, e.g., by region or by age,
  - Computing new variables, e.g., length of days until recovery, or delays in reporting cases
  - Aggregating data, e.g., counting cases or solving rates/means

8

## A Typical Data Science/Analytics Workflow

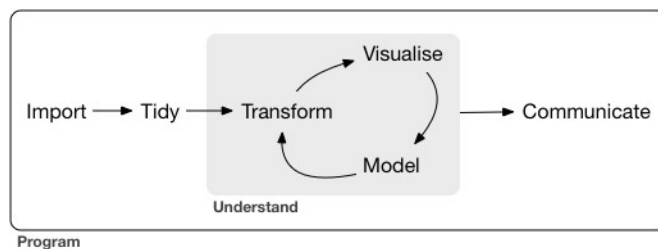


### ► Visualize

- Plot data into graphs so that patterns and features may be explored and insights be extracted from what is seen.

9

## A Typical Data Science/Analytics Workflow

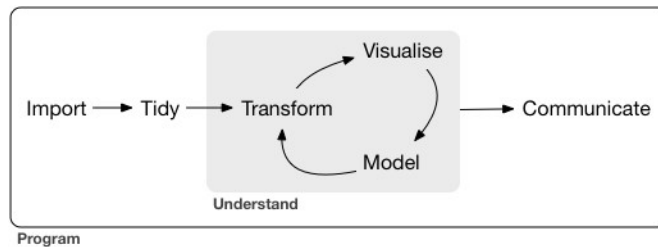


### ► Model

- When necessary, models help in summarizing the complex relationships and the patterns found from visualizations
- Designing models for prediction or forecasting

10

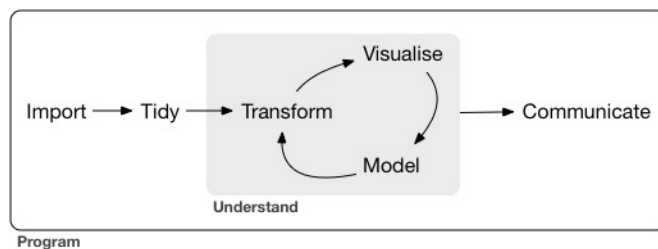
## A Typical Data Science/Analytics Workflow



- Communicate
  - These include:
    - Writing reports, creating dashboards, making presentations, compilations, etc.

11

## A Typical Data Science/Analytics Workflow



- Program
  - All these processes to be encapsulated in a data science project plan
  - Possible to be encapsulated in one software, but it's not impossible to use more than one depending on team members' capabilities to process and analyze data.

12

## Membership to Data-Driven Teams

13



## About the UP Pandemic Response Team

14

## About UP COVID-19 Pandemic Response Team

- **Leaders:**
  - Dr Teodoro Herbosa, UP System
  - Dr Alfredo Mahar Lagmay, Executive Director, UP Resilience Institute
- **Who we are**
  - 200+ experts and volunteers from the whole UP System, from Baguio to Davao
  - spanning multiple fields: political scientists, statisticians, mathematicians, geographers, geologists, medical doctors, linguists, economists, etc.
- Next slides from Dr Lagmay: Our Portfolio



15

# endcov.ph

16

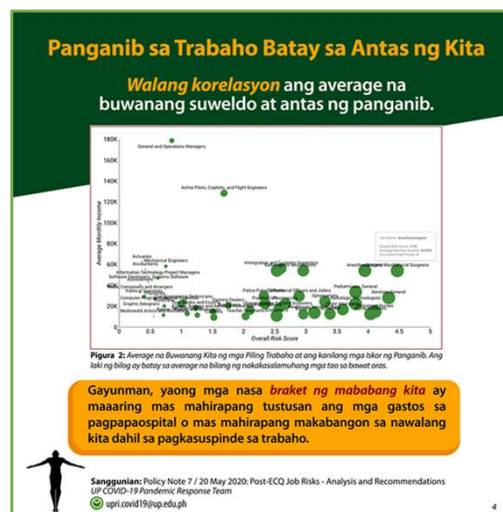
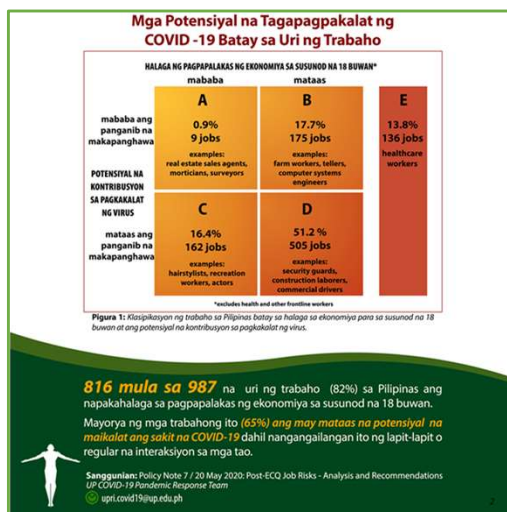


## Yani, the Endcov bot



17

17



To better inform the general public, the team has also redesigned their studies into **animation, video formats and flashcards** to be more accessible.

18

18

UP Resilience Institute  
June 3 at 8:41 PM

Antonaa i maana a "flattening the curve" na miyapantar so curv...  
Inibegay a UP COVID-19 "Manga bababa a Semt Curve".

Batiyaangka saya...  
Published May 8, 2020  
Initogalin saki (Meranaw)  
Noroniah O. Taher  
Musmira C. Bantuas-Ramona  
Miscille T. Domato  
Rayhana P. Ali  
Layout i Ian Villanueva

#UPCOVID19Response

**Antonaa i n... flattening t...**

Giyai so madidiang...  
o isaka inged ka ad...  
kapephakadake...  
COVID-19 para ma...  
ago so mga gumag...  
kagampasi iran ko l...  
miyatangkang a ade...  
iyan a COVID 19.

(Flattening the curve refers measures that keep the da... cases at a manageable level)

Phoon sar: Ang Mga S...

UP Resilience Institute  
May 30 at 9:22 PM

Ano ba iton "epidemic wave"? (Isinalin sa Waray-Leyte)

Ginhahatag iton nga...  
Nahinabato iton Usa...  
Epidemic Wave" han...  
orihinal nga English...  
ngan ginhubad ha Fil...

Basaha an briefer ha...  
Basaha an briefer ha...

**Ano ng...**

Ito ang kur...  
ng bilang r...  
marating

Number of cases

Sa pagbawas...  
bababa ang ku...

Mula sa "Palina...  
Isinalin sa Waray...  
UP COVID-19 Res...  
sign: upcovid19

**Paliwanag Kung Paano Nagaganap ang Isang Epidemya at Ano ang Ibig Sabihin ng Epidemic Wave**

Mary Grace Dacuma, Ph.D.  
University of the Philippines, Los Baños  
(salin sa Filipino ng orihinal na Ingles)

Ipinagpapalagay na nagsimula ang epidemiyang COVID-19 sa Pilipinas nang magkaroon ng imported na kaso mula sa isang tao o mga taong may impekasyon na pumasok sa Pilipinas. Ang tao o mga taong may virus na ito ay ang pinakaunang natukoy na kaso ng nakalahawang sakit (index case) na nagkakat ng virus sa iba dito sa bansa. Hindi kabilang ang (mga) index case, lalo na yung hindi naman naging dahilan ng lokal na transmasyon, sa bugso ng epidemiyang (epidemic wave).

Kapag nagkaroon ng lokal na transmasyon ng virus sa ibang tao, may panahon ng ingkubasyon (ibig sabihin nito na nawala na ng virus ang isang tao ngunit wala pang anumang sintomas o klinikal na mga senyales) na karanawang 5.2 araw hanggang 14 na araw lumalabas. Iyon ang dahilan kung bakit may patag na linya pagkaraang magkaroon ng index case. (tingnan ang Figural).

**Policy notes** are now available in English as well as Tagalog, Ilokano, Bikol Sentral, Waray, Cebuano, Hiligaynon, Aklanon, Kapampangan, Itawis, Chavacano de Zamboanga, Meranaw, and Bahasa Sug.

19

**EPIDEMIC PEAK DYNAMICS**

shifted peak (to buy us more time to prepare)

community quarantine period

reproductive number after the community quarantine period

- R=0.9
- R=1.5
- R=2
- R=2.5
- R=3
- R=4

**Time-Varying Reproduction Number  $R_t$**

Some Results: National  $R_t$  as of 16 April 2020

Philippines  $R_t$

Mar 16 Mar 23 Mar 30 Apr 06 Apr 13

**COVID-19 Threshold Model**

The response team has drafted a number of **policy notes** based on the results of their studies, including recommendations for a graduated activation of the ECQ that depends on the level of risk per area.

20

## Working with the national government

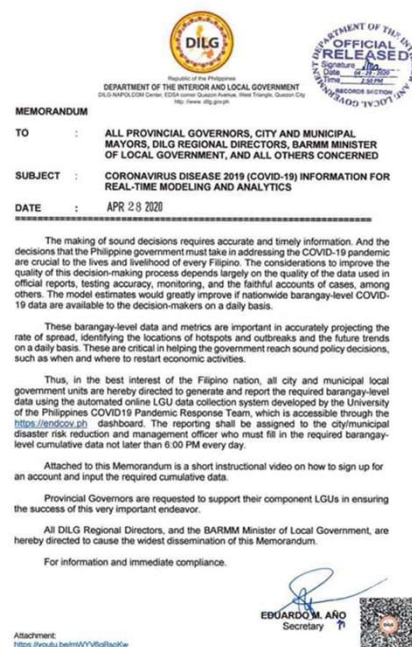
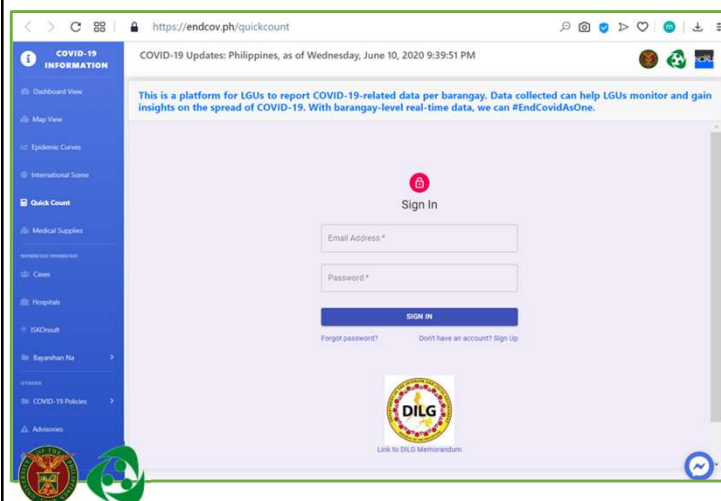
Team members have also **given presentations to the President and the Inter-Agency Task Force on Emerging Infectious Disease (IATF)**. Given their highly specialized competencies, members have been called to join the IATF Technical Working Group on Anticipatory and Forward Planning.



21

21

## Coordination with the LGUs



22

22



Aside from working with national and local government units, the UP response team has been collaborating with academics from the National University of Singapore, University of California Davis and University College London as well as local academic institutions.

On the international front, the team has joined the Forecast-based Warning, Analysis and Response Network (FOREWARN), an organization of academics, scientists and humanitarian workers.



23

23



24

## About UP COVID-19 Pandemic Response Team

### ► Projects I am Involved In:

1. Epidemic Curves of the UP COVID-19 Pandemic Response Team
2. ENDCOV Map on current active cases and probability of outbreak
3. LGU Data Analytics Group
4. Compendium of COVID-19 Statistics for Island Groups, Regions, Provinces, Cities, and Municipalities
5. PSPHP Graphs, Time Varying-R Dashboard with the LEADS 4 Health Security and Resilience



25

## About UP COVID-19 Pandemic Response Team

### ► Platform:

- <https://endcov.up.edu.ph>

### ► Contact Us:

- [upri.covid19@up.edu.ph](mailto:upri.covid19@up.edu.ph)



26



## About the L4H Consortium

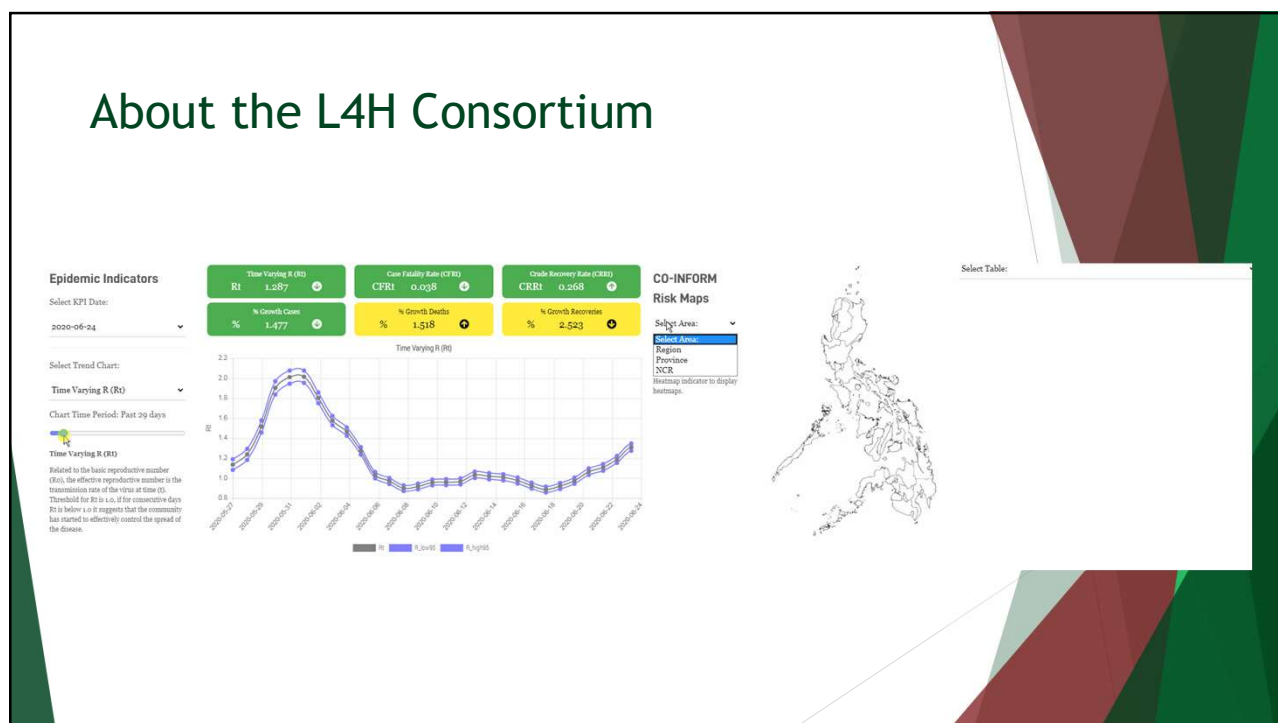
27

## About the L4H Consortium

- ▶ Leading Evidence-based Actions through Data Science for Health Security and Resilience
- ▶ A consortium of data scientists, physicians, mathematicians, and epidemiologists, convened by the Philippine Society of Public Health Physicians (PSPHP)
- ▶ website: <https://leads4health.org/l4h/>

28

## About the L4H Consortium



29

## About the L4H Consortium

### Our Members

- **Jason V. Alacapa, MD, MBA, MPH, MHM** - Consortium Co-convenor, Philippine Society of Public Health Physicians; CEO, metaHealth Insights and Innovation, Inc.
- **Geminn Louis C. Apostol, MD, MBA** - Assistant Professor and Environmental Health Specialist, Ateneo School of Medicine and Public Health
- **DJ Darwin R. Bandy, DVD, PhD(c)** - Assistant Professor, University of the Philippines
- **Peter Julian A. Cayton, PhD** - Associate Professor II, University of the Philippines Diliman
- **Lester Sam Araneta Geroy, MD, MPH, MSc** - President, Philippine Society of Public Health Physicians
- **Dominic Ligot** - Chief Technology Officer, CirroLytix Research Services; Board Member, AAP; Board Member, PCIJ; Lecturer, University of Asia and the Pacific
- **Jason D. Ligot, MD** - Director, Organic Intelligence
- **Robert Neil F. Leong, MSc, PhD(c)** - Chief Data Officer, metaHealth Insights and Innovation, Inc.; PhD Researcher, The University of New South Wales - Sydney; Assistant Professor 2, Mathematics and Statistics Dept. - De La Salle University
- **Lionel A. Peters, MD, MPM** - Member, Philippine Society of Public Health Physicians
- **Michael Angelo B. Promentilla, PhD** - Professor and Research Lead, Waste and Resource Management Unit, Center for Engineering and Sustainable Development Research, De La Salle University
- **Jomar F. Rabajante, PhD** - Professor, University of the Philippines Los Baños and University of the Philippines Open University
- **Miguel Antonio S. Salazar, MD, MSc, Dr.sc.hum.(c)** - Consortium Co-Convenor, Philippine Society of Public Health Physicians
- **Xerxes T. Seposo, MENRM, MPH, PhD** - Assistant Professor, Nagasaki University
- **Jan Gil G. Sarmiento, MSc (ongoing)** - Instructor, University of the Philippines Diliman
- **Theresa Rosario Tan, MABA (ongoing)** - Associate Consultant, CirroLytix Research Services
- **April Anne S. Tigie, MSc, PhD(c)** - Research staff, Waste and Resource Management Unit, Center for Engineering and Sustainable Development Research, De La Salle University

30

## About the L4H Consortium

Our Institutions



31

## Example Projects

32



## Example Projects

1. Epidemiological Statistics for COVID-19 (w/UP PRT)
  - > Compiling epidemiological statistics for website and PDF statistical reports outputs
2. Short-term Forecasting of COVID-19 Cases (w/ L4H)
  - > Forecasting for national and regional reported case counts of COVID-19
3. Mapping COVID-19 Cases (w/ UP PRT)
  - > Mapping the distribution of active cases and case per 100,000 population by barangay in the Philippines

33

## Demonstrating the Workflow with a Team using Example Projects

34

## Demonstrating the Workflow with a Team using Example Projects

Project	Import
Epidemiological Statistics for COVID-19 (w/UP PRT)	Primary Data Source: DOH Data Drop for the Case Information Dataset [ <a href="https://bit.ly/DataDropPH">https://bit.ly/DataDropPH</a> ]  Secondary Data Source:
Short-term Forecasting of COVID-19 Cases (w/ L4H)	PSGC Data with 2020 Population Counts [ <a href="https://psa.gov.ph/classification/psgc/downloads/2_PSGC%20Q-2022-Publication-Datafile.xlsx">https://psa.gov.ph/classification/psgc/downloads/2_PSGC%20Q-2022-Publication-Datafile.xlsx</a> ]
Mapping COVID-19 Cases (w/ UP PRT)	Researchers: Trixie Delmendo, UPRI; Peter Julian Cayton, UP PRT

35

## Demonstrating the Workflow with a Team using Example Projects

The diagram illustrates a workflow for data import. It starts with a Google Drive folder named "READ ME FIRST" containing a file "READ ME FIRST (09.17).pdf". An arrow points from this folder to a screenshot of the Philippine Standard Geographic Code (PSGC) website. Another arrow points from the website to a screenshot of a Google Drive folder containing a list of files, including "DOH Data Drop (09.17).pdf" and "PSGC Data with 2020 Population Counts (09.17).xlsx".

36

## Demonstrating the Workflow with a Team using Example Projects

Project	Tidy
Epidemiological Statistics for COVID-19 (w/UP PRT)	Procedures: <ul style="list-style-type: none"> <li>➤ Reconcile and correct residence variable matched with PSGC data, assuming Barangay PSGC as supreme location variable</li> <li>➤ Reconcile between dataset versions to add missing recovery dates</li> </ul>
Short-term Forecasting of COVID-19 Cases (w/ L4H)	Researchers: Trixie Delmendo, UPRI; Jan Gil Sarmiento, UPD SS Peter Julian Cayton, UP PRT
Mapping COVID-19 Cases (w/ UP PRT)	

37

## Demonstrating the Workflow with a Team using Example Projects

```

54 },"Provinces"] <- "Cotabato City"
55 doh_data[which(!is.na(doh_data$Provinces) &
56 doh_data$Provinces=="Samar (Western Samar)","Provinces"] <-
57 "Samar"
58 doh_data[which(!is.na(doh_data$Provinces) &
59 doh_data$Provinces=="Cotabato (North Cotabato)","Provinces"]
60 <- "Cotabato"
61 # unique(doh_data$Provinces) %>% unique(places$Admin2_Name)
62 ## Fixing the Provinces for NCR cases with NA in Provinces
63 doh_data[which(is.na(doh_data$Provinces) & doh_data$RegionRes=="NCR"
64 ],"Provinces"] <- "METRO MANILA"
65 # Replacing Negros Island region -> Region 6/7
66 places[which(places$Admin2_Name=="Negros Occidental"),3] <- "Region
67 VI"
68 places[places$Admin2_Name=="Negros Oriental",3] <- "Region VII"

```

R version 4.1.0 (2021-05-18) -- "Camp Pontanezen"  
 Copyright (C) 2021 The R Foundation for Statistical Computing  
 Platform: x86\_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
 You are welcome to redistribute it under certain conditions.  
 Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.  
 Type 'contributors()' for more information and  
 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
 'help.start()' for an HTML browser interface to help.  
 Type 'q()' to quit R.

38

## Demonstrating the Workflow with a Team using Example Projects

Project	Transform
Epidemiological Statistics for COVID-19 (w/UP PRT)	<p>Procedures:</p> <ul style="list-style-type: none"> <li>➤ Group data by geographic hierarchy: 1) National, 2) Island Group, 3) Regional, 4) Provincial, 5) City/Municipality</li> <li>➤ Compute for cumulative and new counts for cases, recoveries, deaths, and active cases per date per hierarchy</li> <li>➤ Compute for case rates such as recovery rates, fatality rates, and doubling times</li> </ul> <p>Researchers: Jan Gil Sarmiento, UPD SS; Peter Julian Cayton, UP PRT</p>
Short-term Forecasting of COVID-19 Cases (w/ L4H)	<p>Procedures:</p> <ul style="list-style-type: none"> <li>➤ Group data by geographic hierarchy: 1) National and 2) Regional</li> <li>➤ Compute for new counts for cases per date per hierarchy</li> </ul> <p>Researchers: Jan Gil Sarmiento, Simon Bismonte, Maryliz Zubiri, and Nicole Uy</p>
Mapping COVID-19 Cases (w/ UP PRT)	<p>Procedures:</p> <ul style="list-style-type: none"> <li>➤ Group data by geographic location: Barangay via PSGC</li> <li>➤ Compute for total counts for cases, recoveries, deaths, and active cases per barangay</li> <li>➤ Compute for case per 100,000 population using PSGC population data and case counts by barangay.</li> </ul> <p>Researchers: Trixie Delmendo, UP RI; Peter Julian Cayton, UP PRT</p>

39

## Demonstrating the Workflow with a Team using Example Projects

```

193
194 #function to count number of confirmed cases ----
195 count.conf <- function(df, sd=NULL, ed){
196   confirmed <- data.frame(table(df$dateRepConf))
197   if(dim(confirmed)[1]==0){
198     return(confirmed)
199   } else {
200     colnames(confirmed) <- c("Date", "Freq")
201     confirmed$date <- as.Date(confirmed$date, date_format = "%Y-%m-%d")
202     pad.confirmed <- pad(confirmed, interval="day", start_val = sd,
203                          end_val = ed)
204     pad.confirmed <- fill_by_value(pad.confirmed, value = 0)
205     return(pad.confirmed)
206   }
207 }
208

```

R version 4.1.0 (2021-05-18) -- "Camp Pontanzen"  
Copyright (c) 2021 The R Foundation for Statistical Computing  
Platform: x86\_64-w64-mingw32/x64 (64-bit)  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

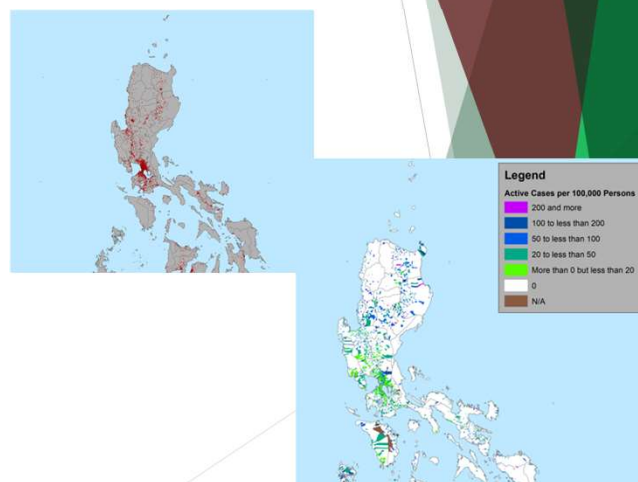
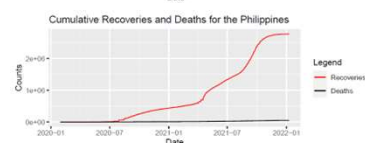
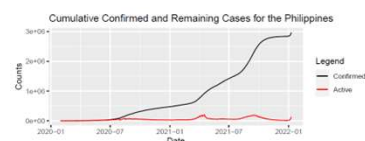
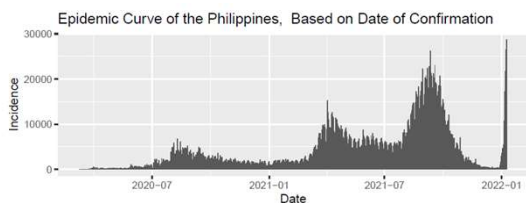
40

## Demonstrating the Workflow with a Team using Example Projects

Project	Visualize
Epidemiological Statistics for COVID-19 (w/UP PRT)	Procedures: ➤ Generate plots of cases, recoveries, deaths, and active cases per hierarchy  Researchers: Jan Gil Sarmiento, UPD SS; Peter Julian Cayton, UP PRT
Short-term Forecasting of COVID-19 Cases (w/ L4H)	Procedures: ➤ Generate plots of cases per hierarchy  Researchers: Jan Gil Sarmiento, Simon Bismonte, Maryliz Zubiri, and Nicole Uy
Mapping COVID-19 Cases (w/ UP PRT)	Procedures: ➤ Produce a map of active case distribution in which case dots are randomly distributed within their barangay of residence ➤ Produce a color map of case per 100,000 by barangay of residence with the barangay colors done by belonging into an interval  Researchers: Trixie Delmendo, UP RI; Steffanie Chua, UP RI; Jake Mendoza, UP RI; Feye Andal, UP RI; UP Resilience Youth Mappers; Peter Julian Cayton, UP PRT

41

## Demonstrating the Workflow with a Team using Example Projects



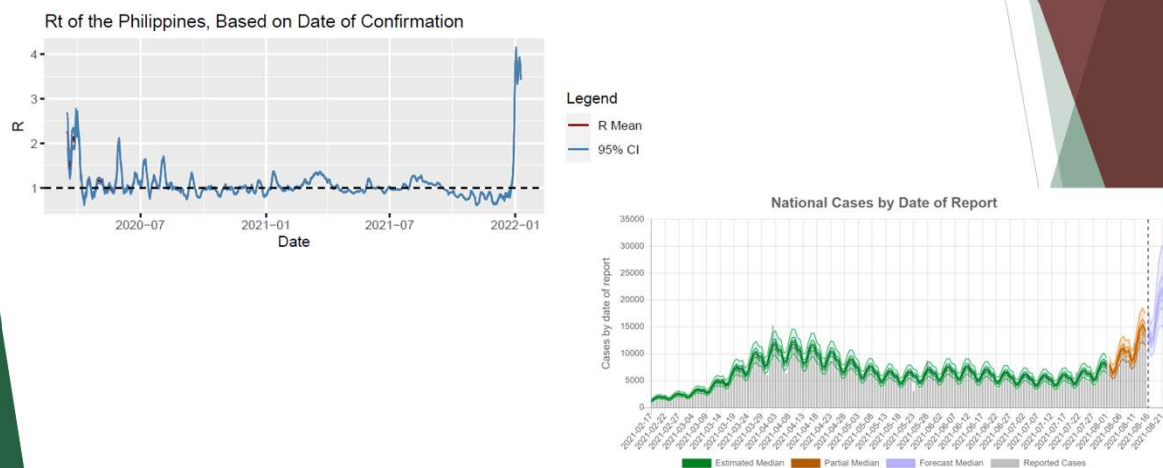
42

## Demonstrating the Workflow with a Team using Example Projects

Project	Model
Epidemiological Statistics for COVID-19 (w/UP PRT)	Procedures: ➤ Generate the time-varying reproduction number based on an existing Bayesian estimation technique  Researchers: Jan Gil Sarmiento, UPD SS; Peter Julian Cayton, UP PRT
Short-term Forecasting of COVID-19 Cases (w/ L4H)	Procedures: ➤ Generate the short-term forecasts assuming a model based on the Epiforecasts methodology  Researchers: Jan Gil Sarmiento, Simon Bismonte, Maryliz Zubiri, Peter Julian Cayton, Robert Neil Leong
Mapping COVID-19 Cases (w/ UP PRT)	(none)

43

## Demonstrating the Workflow with a Team using Example Projects



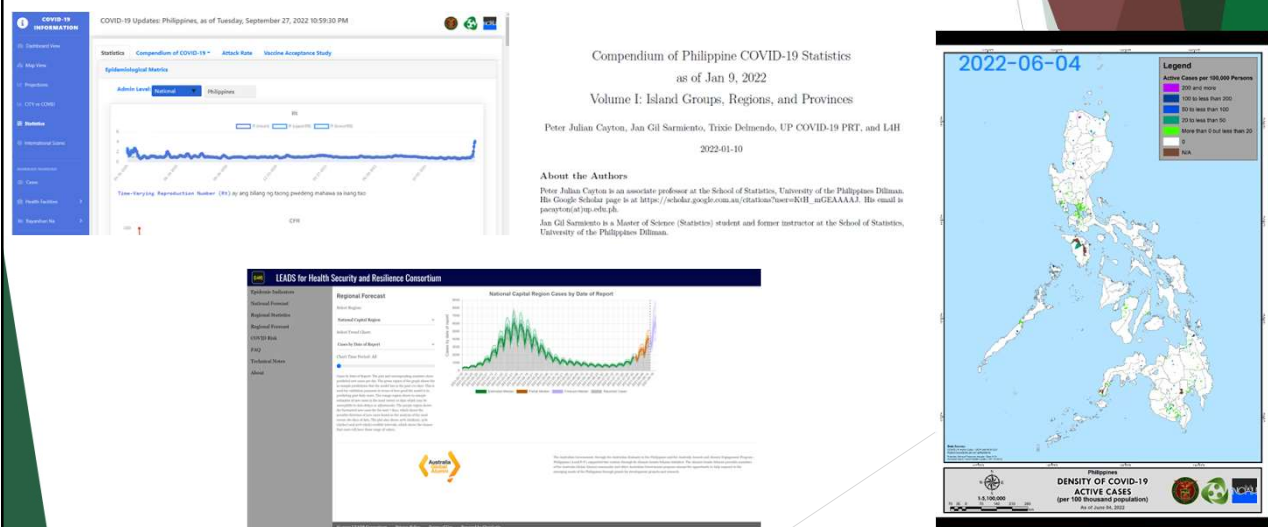
44

## Demonstrating the Workflow with a Team using Example Projects

Project	Communicate
Epidemiological Statistics for COVID-19 (w/UP PRT)	<p>Procedures:</p> <ul style="list-style-type: none"> <li>➢ Publish a website page for uploading the statistics</li> <li>➢ Generate a PDF report file</li> </ul> <p>Researchers: Feye Andal and developers from the UP RI; Peter Julian Cayton, UP PRT</p>
Short-term Forecasting of COVID-19 Cases (w/ L4H)	<p>Procedure:</p> <ul style="list-style-type: none"> <li>➢ Publish the results through the L4H website</li> </ul> <p>Researchers: Dominic Ligot, Mark Toledo, and Angelica Mhay Salazar</p>
Mapping COVID-19 Cases (w/ UP PRT)	<p>Procedure:</p> <ul style="list-style-type: none"> <li>➢ Produce maps based on the results</li> </ul> <p>Researchers: Trixie Delmendo, UP RI; Steffanie Chua, UP RI; Jake Mendoza, UP RI; mappers and research of the UP RI and UP RI Youthmappers,</p>

45

## Demonstrating the Workflow with a Team using Example Projects



46



## Closing Remarks

47

## Closing Remarks

- ▶ Data science involves not only technical skills such as programming, statistics, and subject matter knowledge, but also communication and team collaboration
- ▶ By making data science projects in teams, the burden per person is reduced and high-quality data products are produced.
- ▶ Utilize each member's comparative advantage skills to maximize quality of output
- ▶ Build cohesion and empathy as a team of analysts and scientists producing data products

48



