# Leveraging R and ChatGPT for Epidemiological Analysis

Prof Dr Jamalludin Ab Rahman MD MPH FPHMM
*Dean, Kulliyyah of Medicine, IIUM*
*President, PPPKAM*

# What is ChatGPT

2



- Developed by OpenAI
- Built on the GPT (Generative Pre-trained Transformer) – Vaswani et al. 2017
- Transformer – a new neural network architecture
- Attention is a mechanism that allows neural networks to focus on specific parts of their input.

# All start with prediction

3

- *I'm going to the......*
- Use the Language Model to predict the next word
- It is already in the smartphone – predictive text
- The language model works on a certain computational framework (neural network, many types, one of them is the Transformer model)

# Example of the 'evolution' in diabetes epidemiology

## Predictive Statistics

**Estimating diabetes risk** using logistic regression based on traditional risk factors (e.g., age, BMI).

## Machine Learning

**Predicting diabetes onset and complications** using algorithms like neural networks and analysing a broader data set.

## Artificial Intelligence

**AI-powered health assistants provide** real-time monitoring, dietary and exercise recommendations, and predict blood sugar levels.

## Generative AI

**Generating synthetic datasets of diabetic patient** records for research using Generative Adversarial Networks (GANs).

# How ChatGPT works in statistics

- Trained from all information available (easier & faster literature search)
- Propose the best* solution, best* method
- Describe concepts faster and easier
- Propose coding/algorithm

* Depending on the trained data

# Prompting tips

1. Persona/role – a lecturer, a professor, a student, a PhD candidate

2. The task/instruction – to improve, to describe, to analyse, to compare etc

3. Expectation/end goals – simple, complex, layman's terms

4. Filter – narrow the output

5. Format the output – table, diagram



ChatGPT Prompting Cheat Sheet — 5 frameworks to level up your prompts. Created by Moritz Kremb (@moritzkremb, thepromptwarrior.com)

ChatGPT is a tool

# What is R

- Open-source programming language environment
- Used especially for **statistical** computing and **graphics**
- Free
- *Steep* learning curve – *this is where ChatGPT is useful*
- Progressive – constant update – *again, this is where ChatGPT is useful*

# How do you rate your skill using R?

ⓘ Start presenting to display the poll results on this slide.

# For this workshop

1.  Good to have two monitors

2.  Install R, download from https://cran.r-project.org

3.  Install RStudio from https://rstudio.com

4.  Learning materials at https://github.com/profjamal/chatgpt

5.  Data from https://github.com/MoH-Malaysia/covid19-public

# The steps

1. Set your objective very clear

2. Understand the data (if they are not yours)

3. Prepare the data – download (link), clean (for missing values, outliers, etc), and visualise the data (in table)

4. Generate the coding from ChatGPT by using proper **prompting**

5. Run the coding in R

6. Verify the results – need your understanding of epid & stat

# Set your analysis objectives

- Begin with the end in mind

- Set clear objectives

- For this workshop, our analysis is on Malaysia's COVID-19 death (using the line listing):

  https://github.com/MoH-Malaysia/covid19-public/blob/main/epidemic/linelist/linelist_deaths.csv

# The prompt, example

| Role | Who are you? | Epidemiologist that is responsible for managing outbreaks in a population of 36,000,000 people |
|---|---|---|
| Task | What do you want? | Generally the aim is to evaluate the deaths from the COVID-19 outbreak in Malaysia. Specifically,<br><br>1) the overall incidence by year<br>2) its distribution based on age, sex, vaccination status, and type of vaccines. |
| Format | What output format do you want? | In R script and/or the visualisations |

# Tips

- Faster if we upload the data to ChatGPT
- Good to proceed in stages
  1. Preview data
  2. Clean data
  3. Analyse base on objective
  4. Request for visualisation
  5. Can even request the narrative
- At each stage, verify the codes and the result of the analysis
- Let the ChatGPT run the whole analysis first (in Phyton)
- Once you are satisfied, then you can ask the code in R, verify again

# Practical

Let's do this together

# Let ChatGPT do all the analysis first

- [https://chatgpt.com](https://chatgpt.com)
- GPT-4, GPT-4o mini, GPT-4o
- Let's start with your first prompt

# Objectives

1. Calculate the overall mortality rate

2. Compare mortality rate by sex, vaccination status, vaccine doses