**Big Data Fundamentals and Applications**

# Statistical Analysis (II)
# Descriptive Statistics – Graph

## Asst. Prof. Chan, Chun-Hsiang

*Master program in Intelligent Computing and Big Data, Chung Yuan Christian University, Taoyuan, Taiwan*
*Undergraduate program in Intelligent Computing and Big Data, Chung Yuan Christian University, Taoyuan, Taiwan*
*Undergraduate program in Applied Artificial Intelligence, Chung Yuan Christian University, Taoyuan, Taiwan*
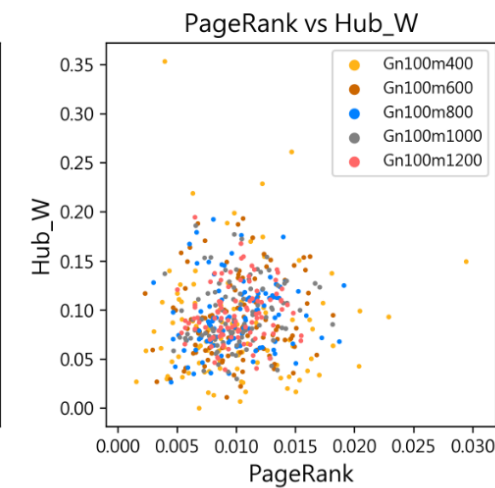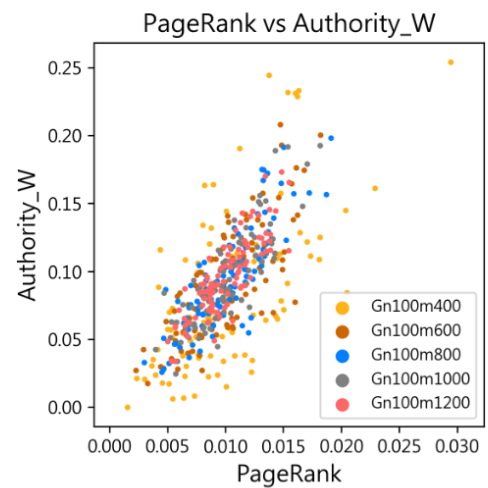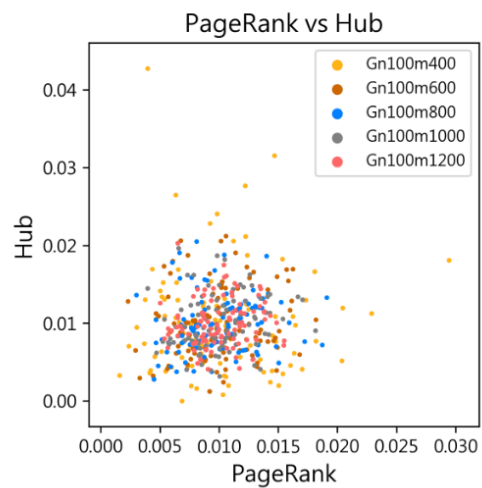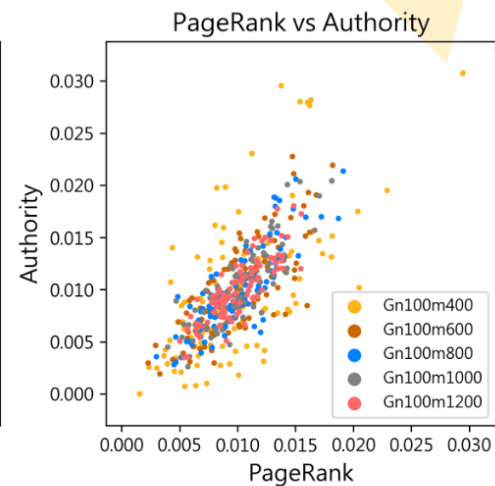
# Outlines

# **Visualization**

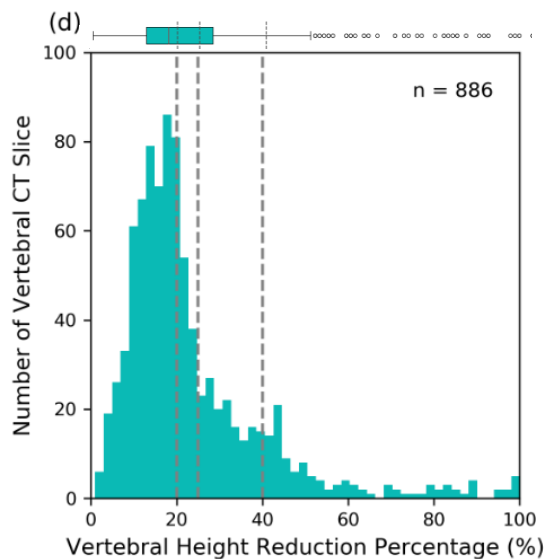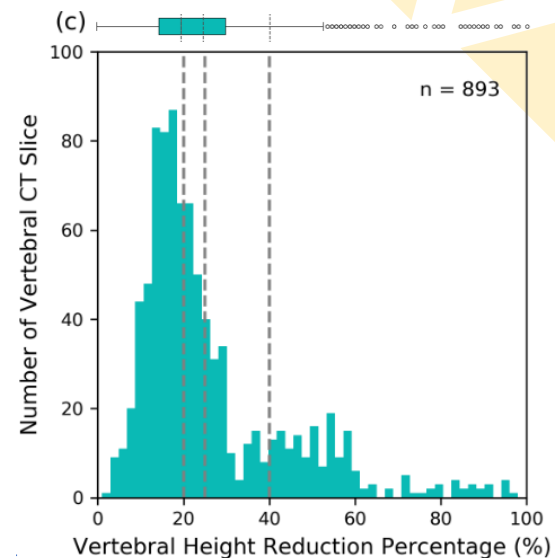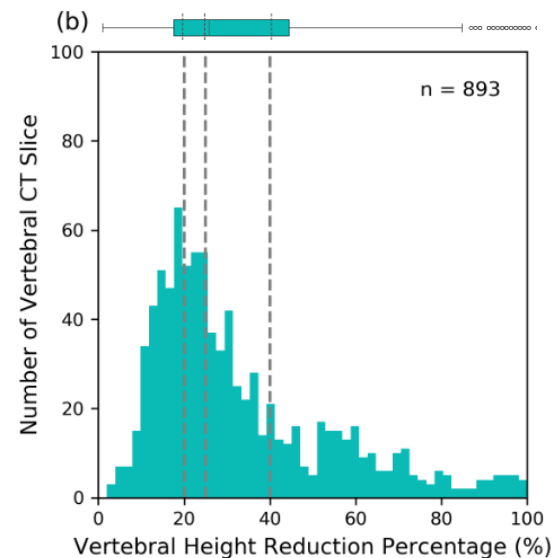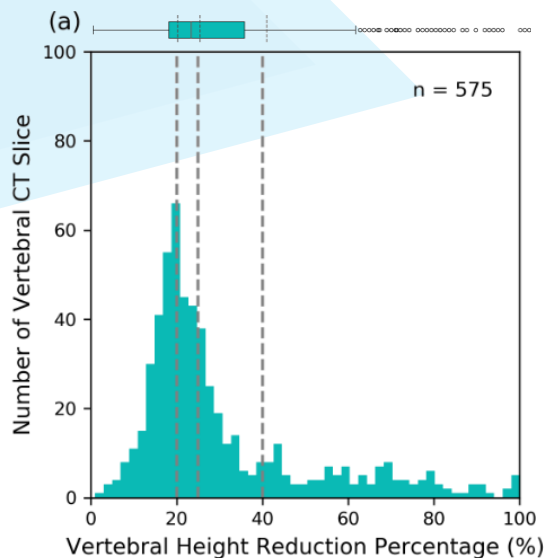- As we mentioned before, data visualization is one of the most effective approach for data exploration to understand the characteristics of each feature or dataset.

- According to the attribute of data, we need to select the appropriate visualization methods for illustration.

- Here, we are going to introduce bar chart, biplot, box plot, control chart, forest plot, histogram, pie chart, Q-Q plot, scatter plot, stem plot, violin plot, radar plot, and line plot.
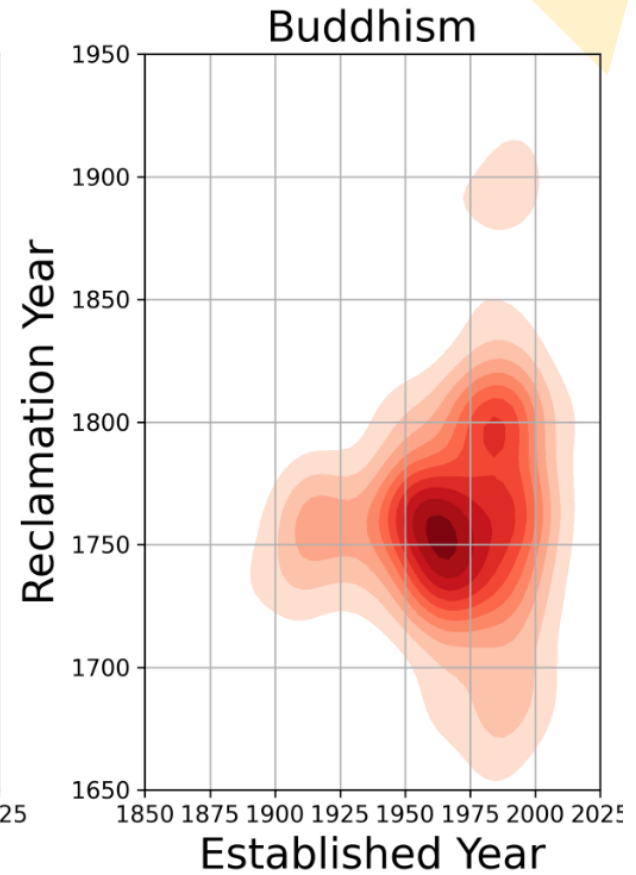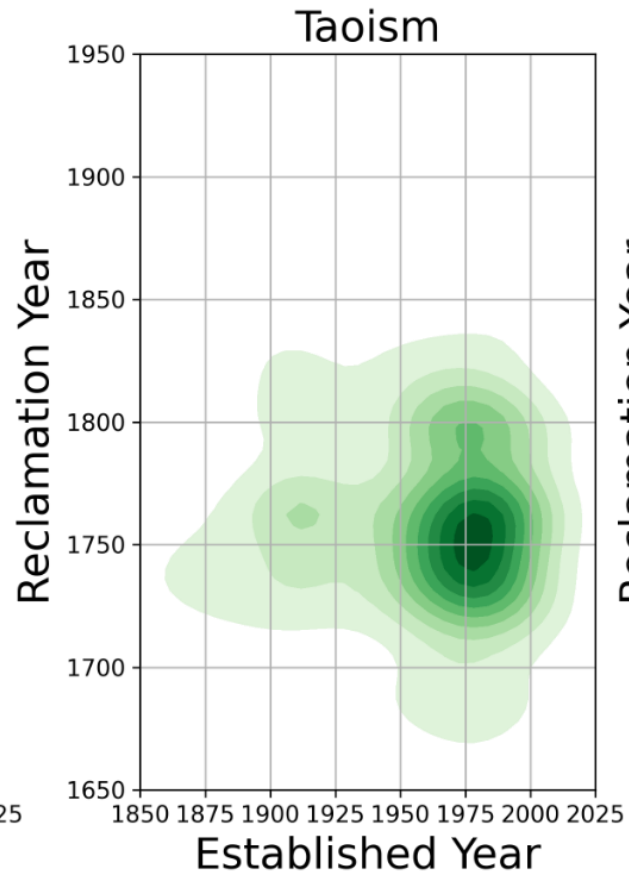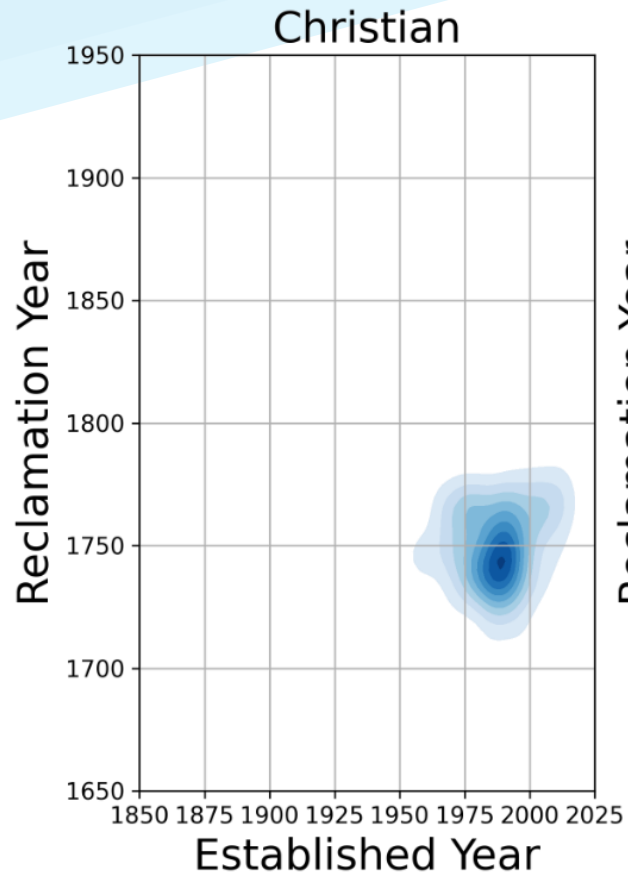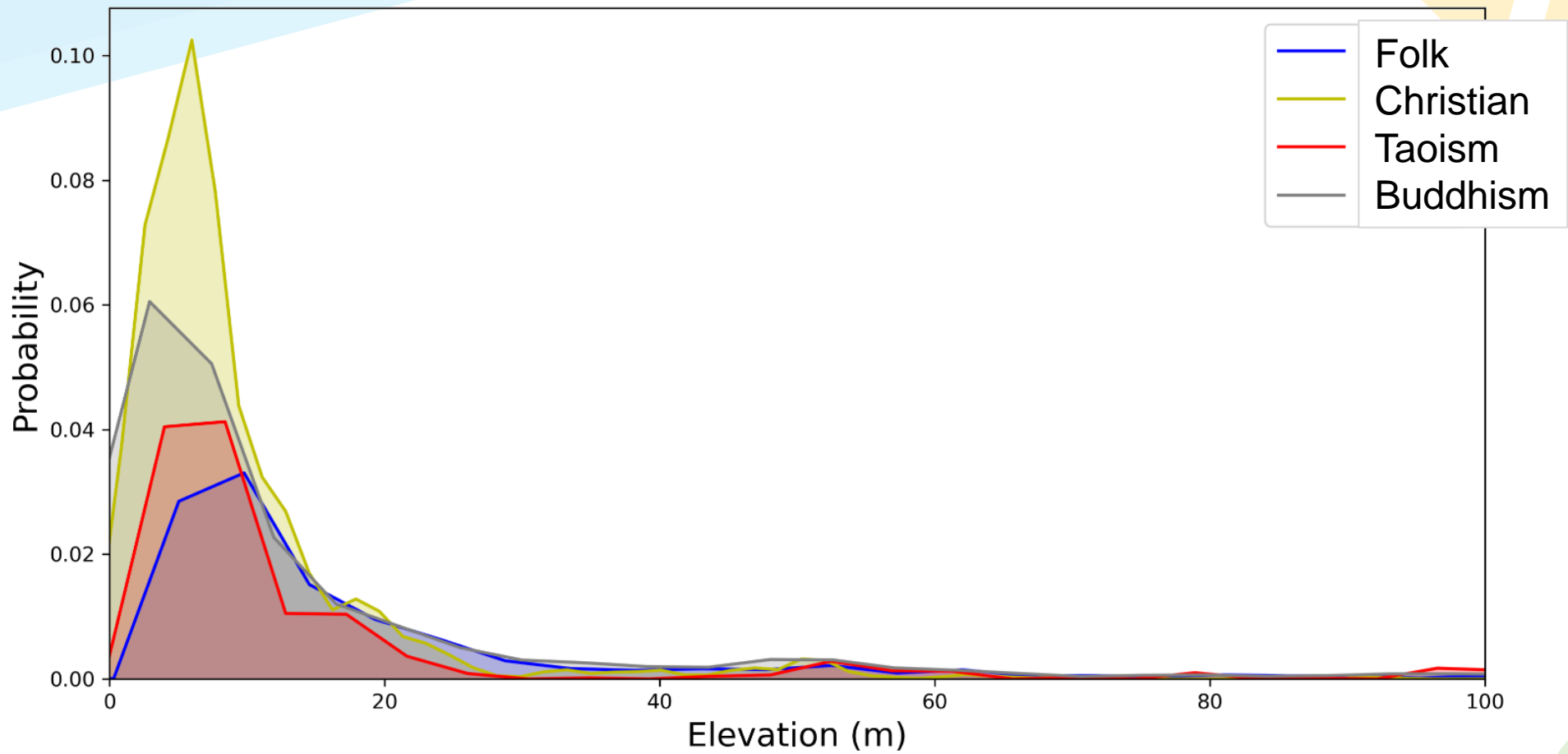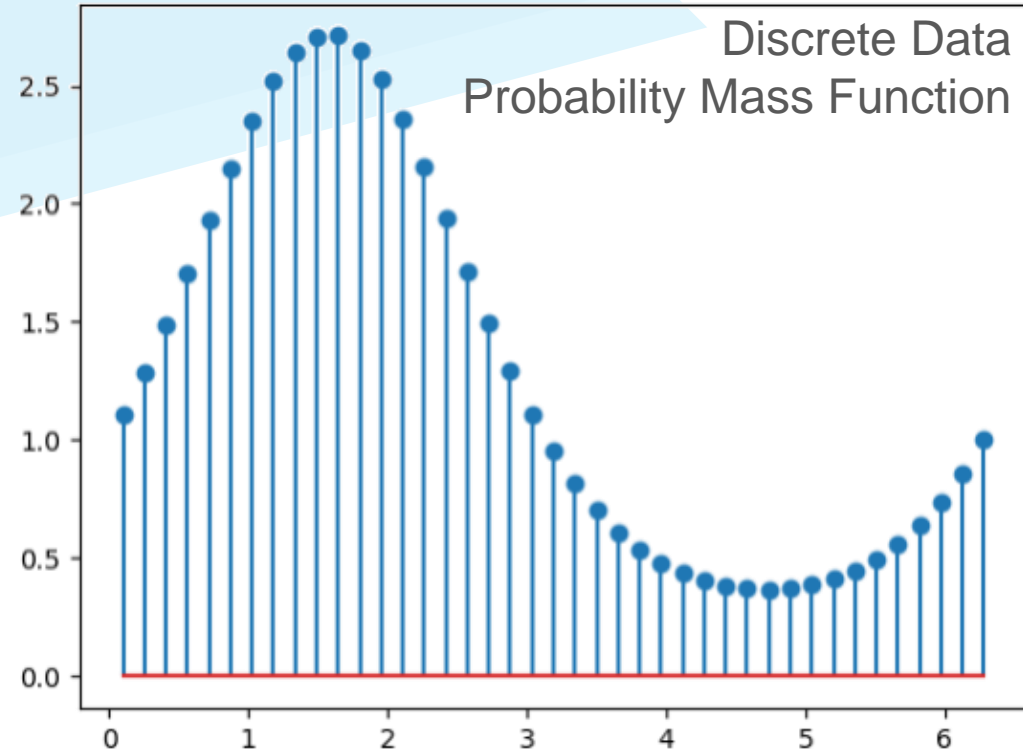
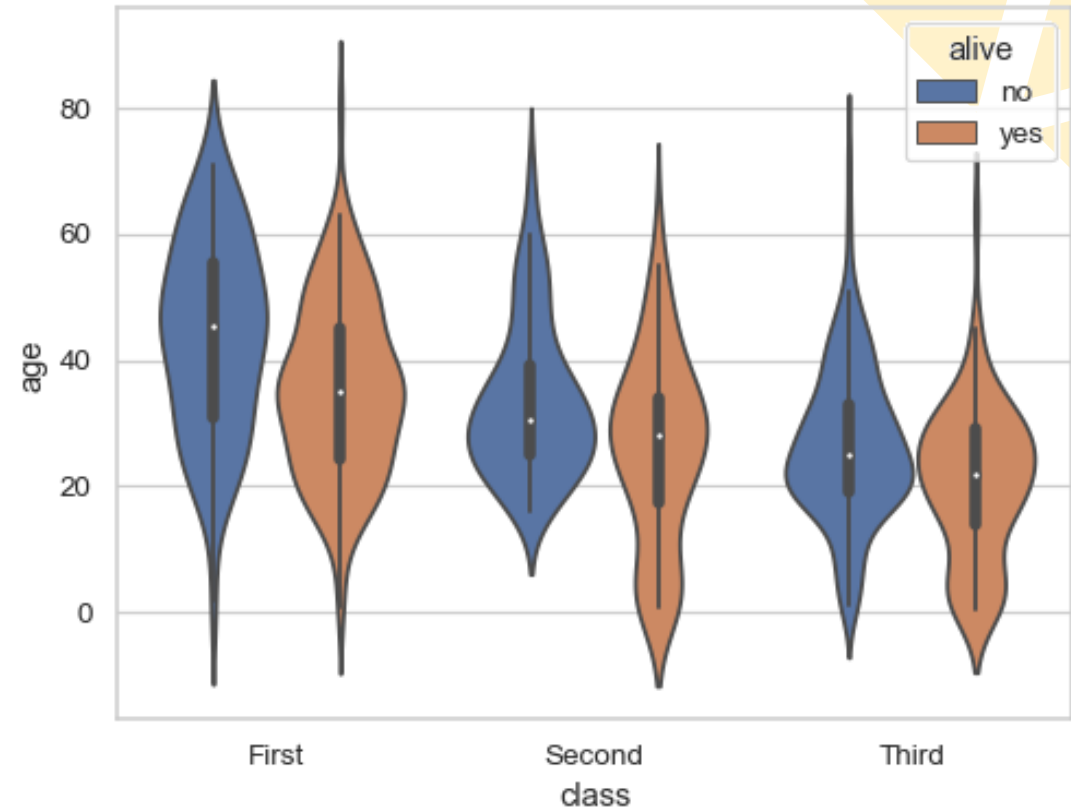# Line Plot & Scatter Plot

# 1D Histogram

# 2D Histogram

# Area Plot

# Stem Plot & Violin Plot



Discrete Data
Probability Mass Function

Continuous Data

**Source:**
https://matplotlib.org/stable/gallery/lines_bars_and_markers/stem_plot.html#sphx-glr-gallery-lines-bars-and-markers-stem-plot-py
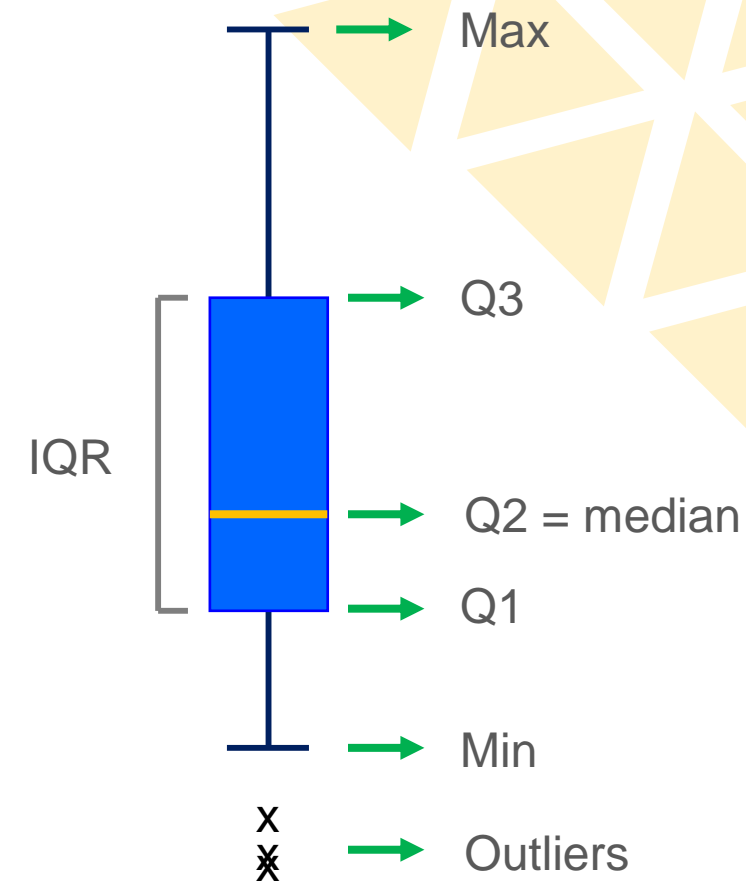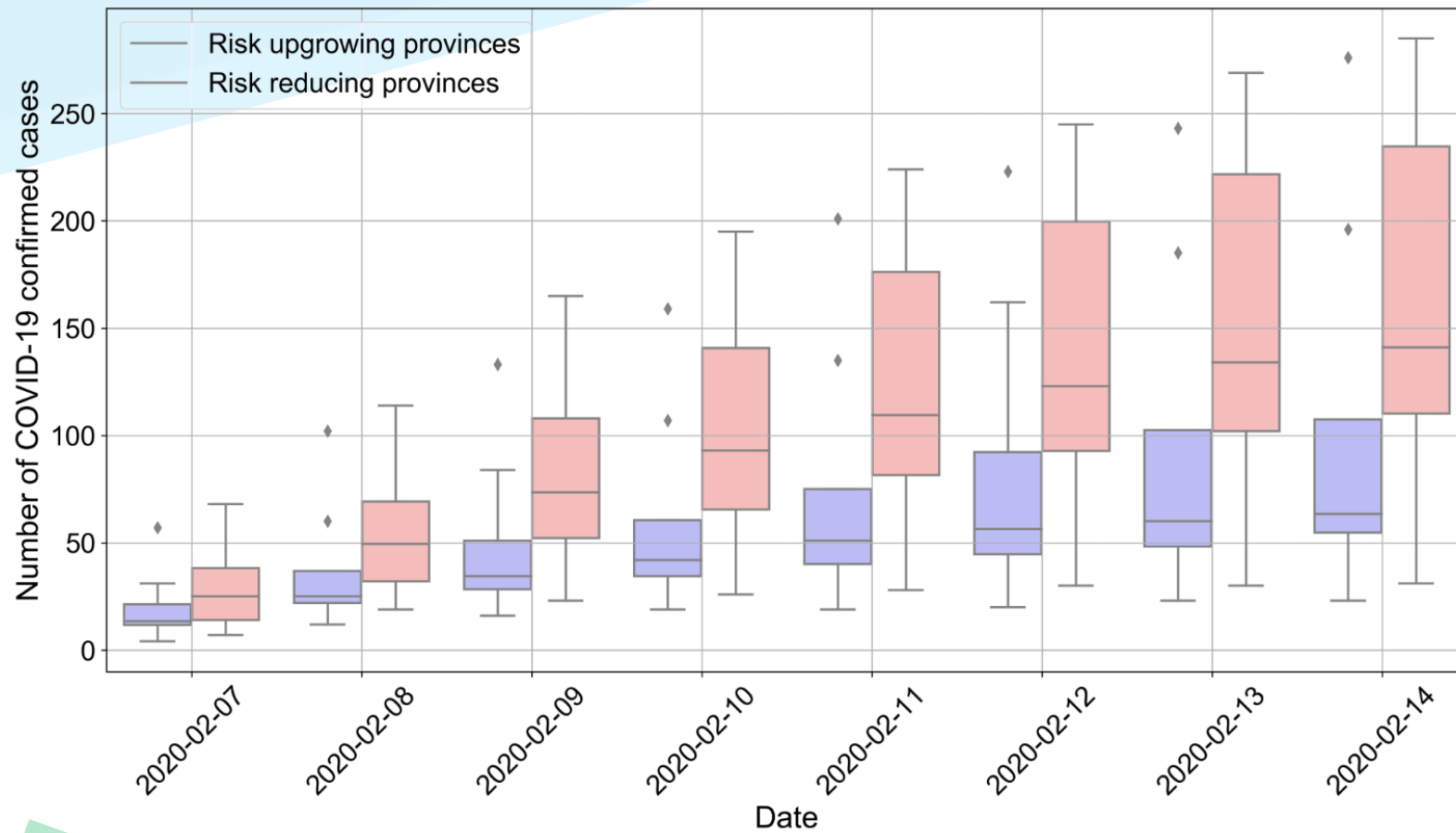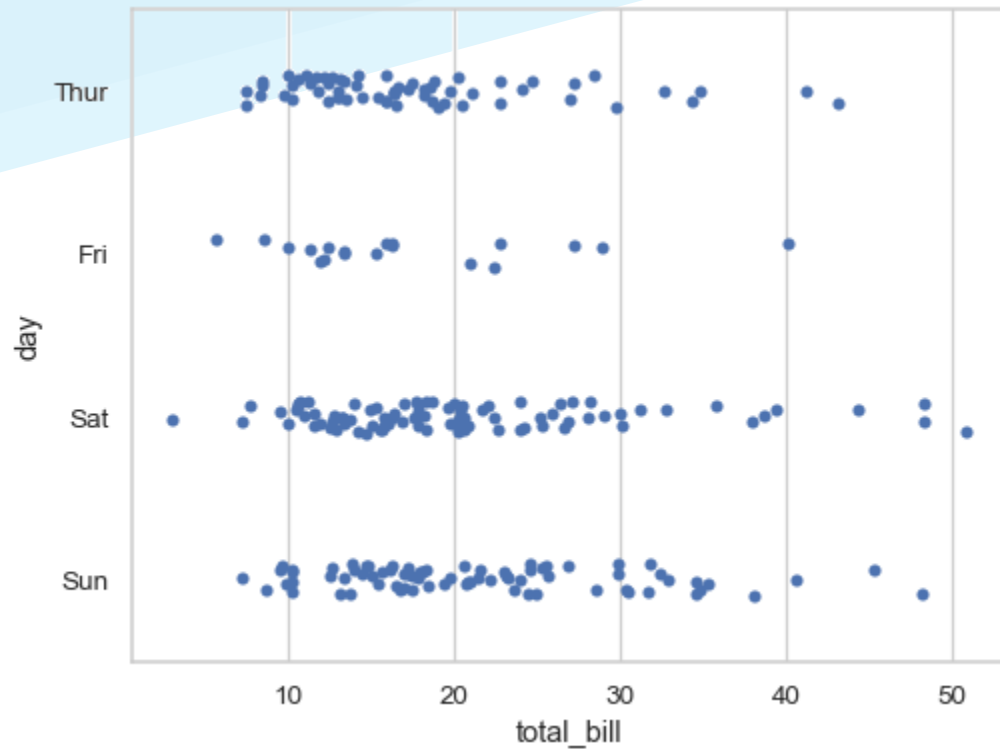
**Source:** https://seaborn.pydata.org/generated/seaborn.violinplot.html

# Box Plot

# Strip Plot & Swarm Plot

Observe the d

# Error Bar Plot

# Pie Chart

# Bar Chart

# Rose Plot & Radar Plot



GOS-1  GOS-2  GOS-3

GOS-4  GOS-5  GOS-6

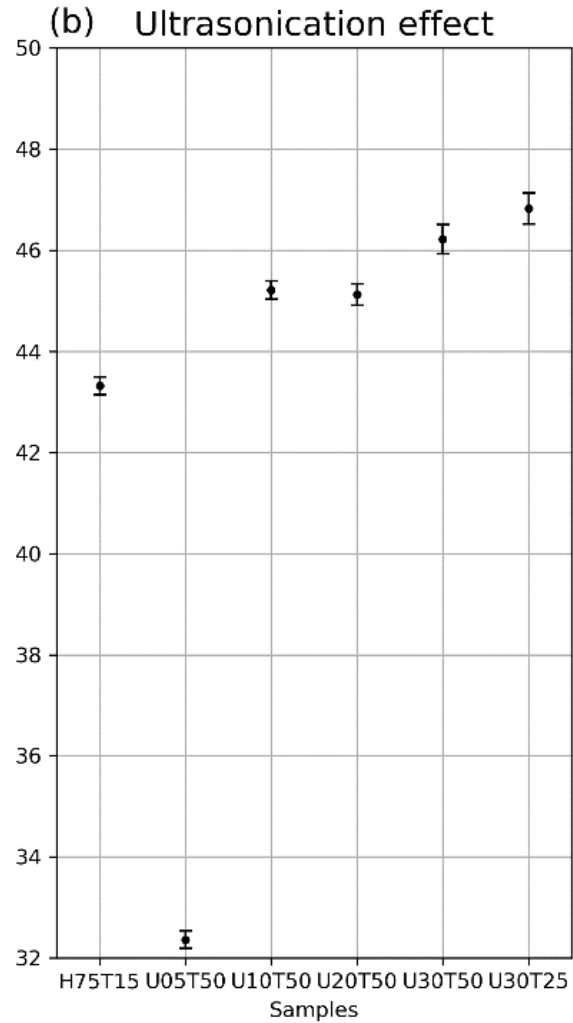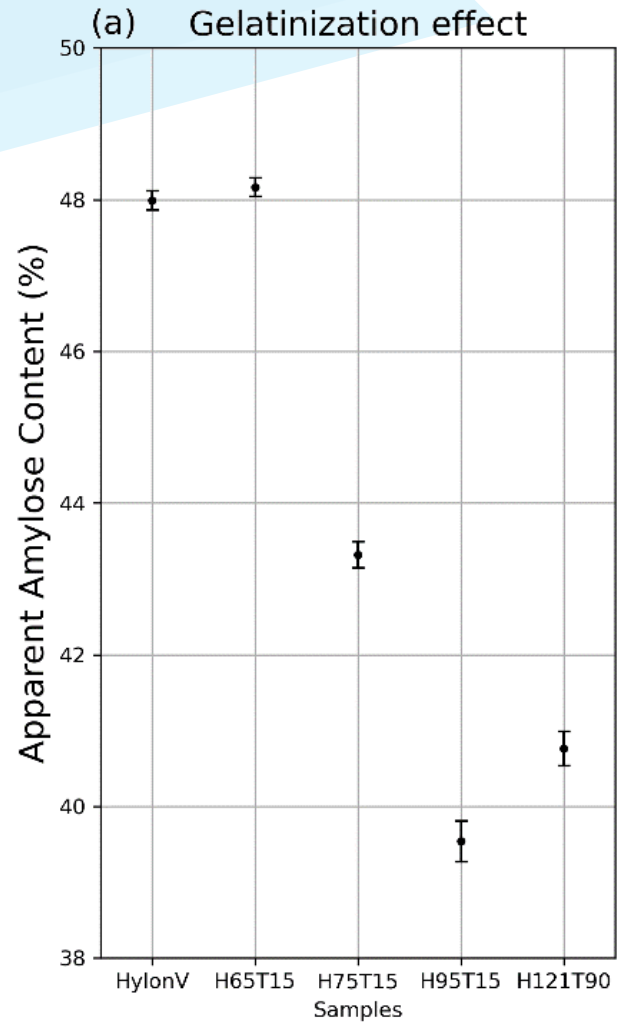**Profile diversity in various GOS samples**

Lin et al. (2022) Profile diversity of galacto-oligosaccharides from disaccharides to hexasaccharides by porous graphitic carbon liquid chromatography-orbitrap tandem mass spectrometry. Food Chem. Vol. 390. 133151.
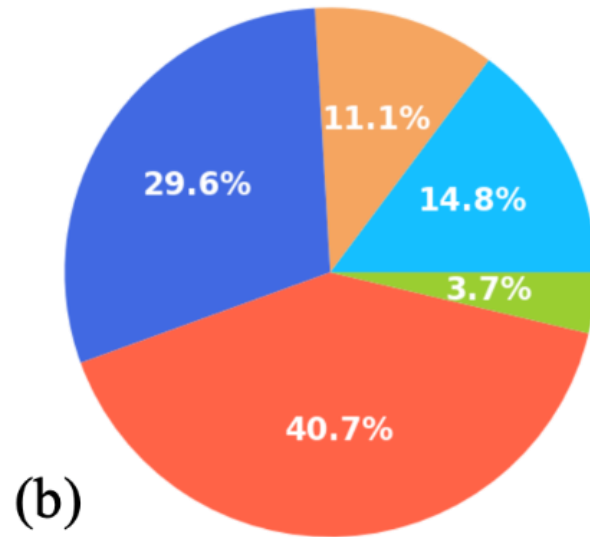


Star Plot of MER IDD and Automated Designs

MER IDD
Design 1
Design 2
Design 3
Design 4

Accuracy
Collision
Link Deflection
Trajectory Completion
Actuator Saturation
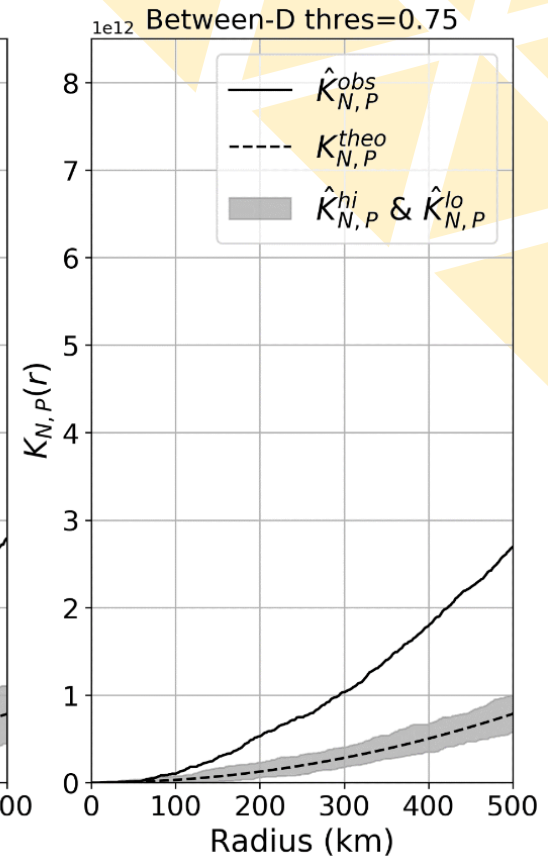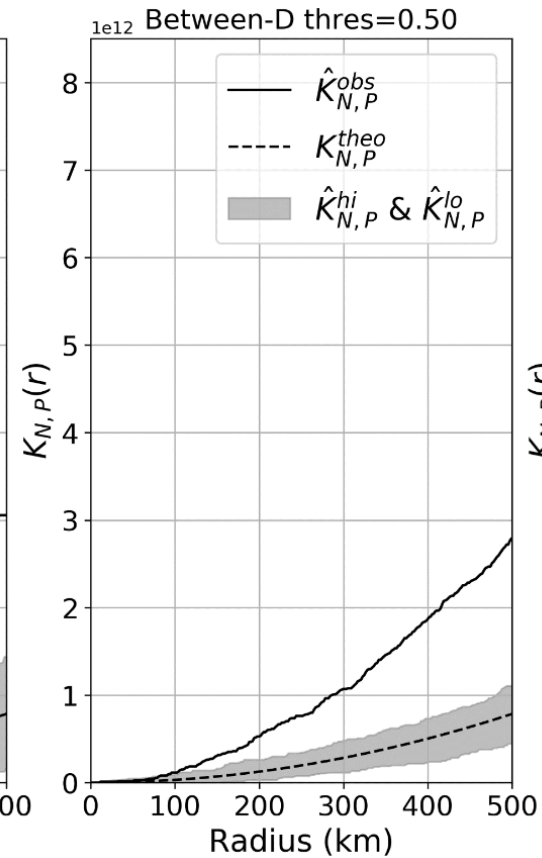Time (Task Completion)
Mass

**Source:** https://en.wikipedia.org/wiki/Radar_chart#/media/File:MER_Star_Plot.gif

# Biplot & Control Chart (I)

# Control Chart (II)

# Forest Plot



| | Diseased | Healthy |
|---|---|---|
| **Exposed** | 20 | 380 |
| **Not Exposed** | 10 | 490 |

$$risk\ of\ developing\ the\ disease\ given\ exposure = \frac{DE}{VE} = \frac{20}{400}$$

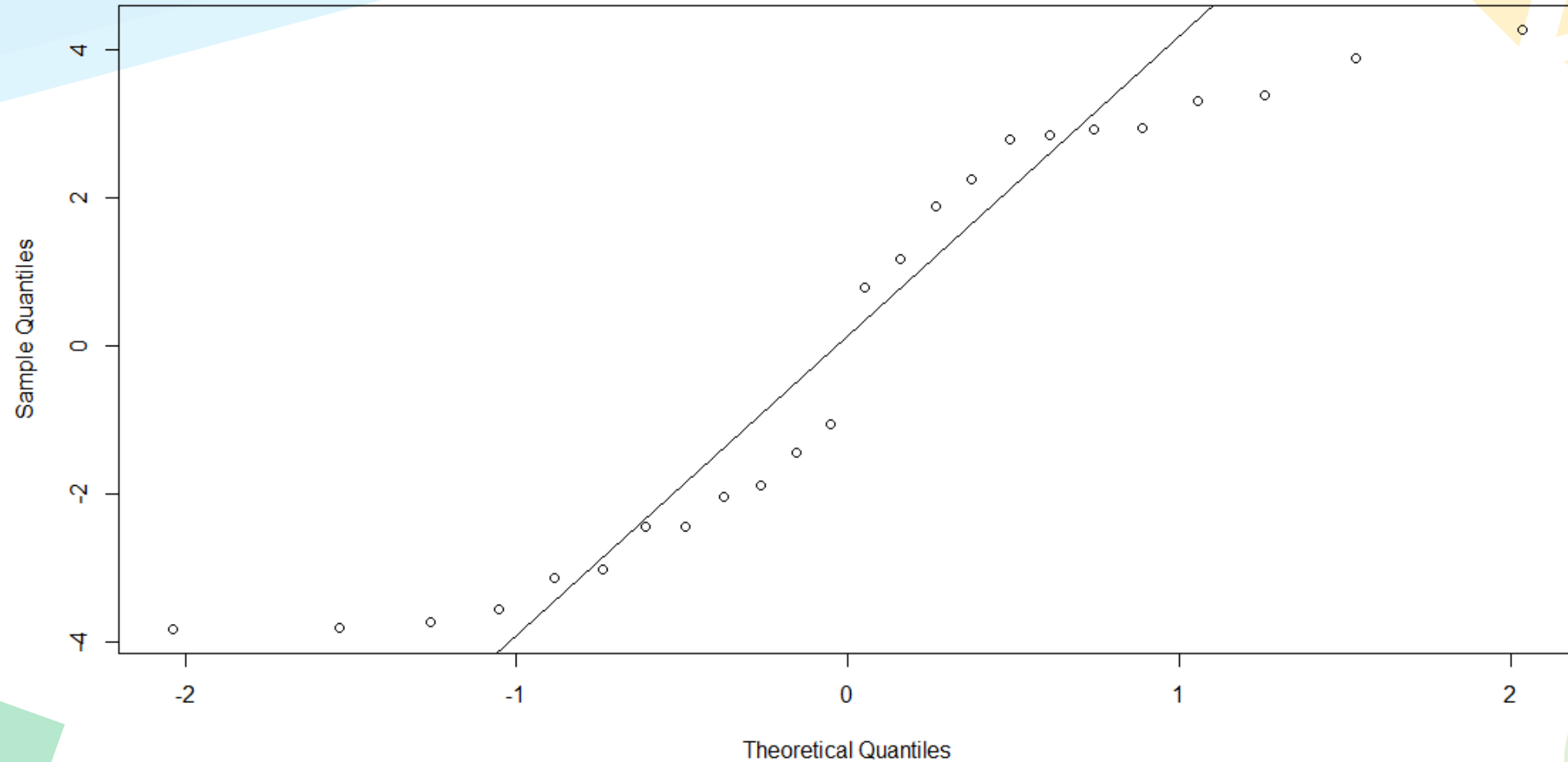$$risk\ of\ developing\ the\ disease\ given\ non-exposure = \frac{DN}{VN} = \frac{10}{500}$$

$$relative\ risk = \frac{\dfrac{DE}{(DE+HE)}}{\dfrac{DN}{(DN+HN)}} = \frac{DE/VE}{DN/VN} = \frac{20/400}{10/500}$$

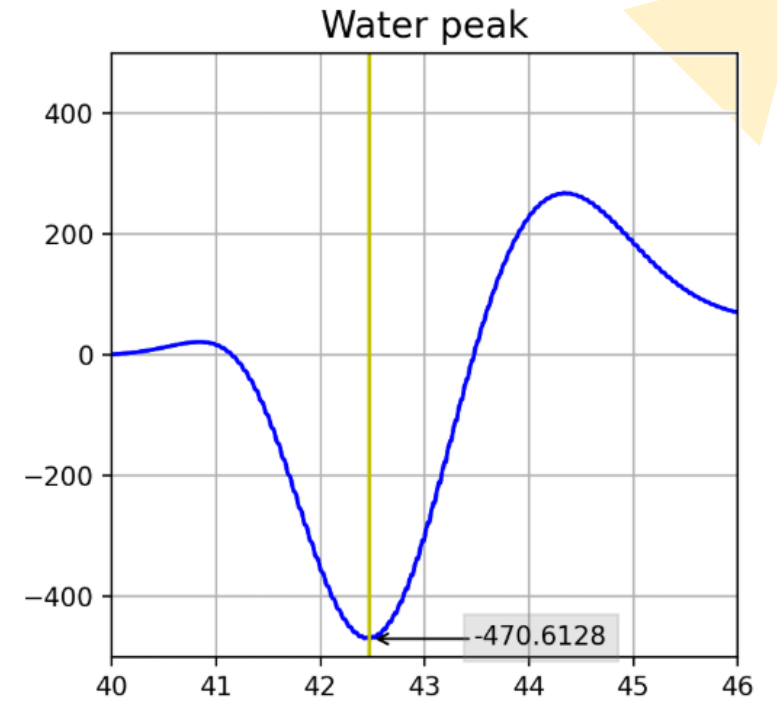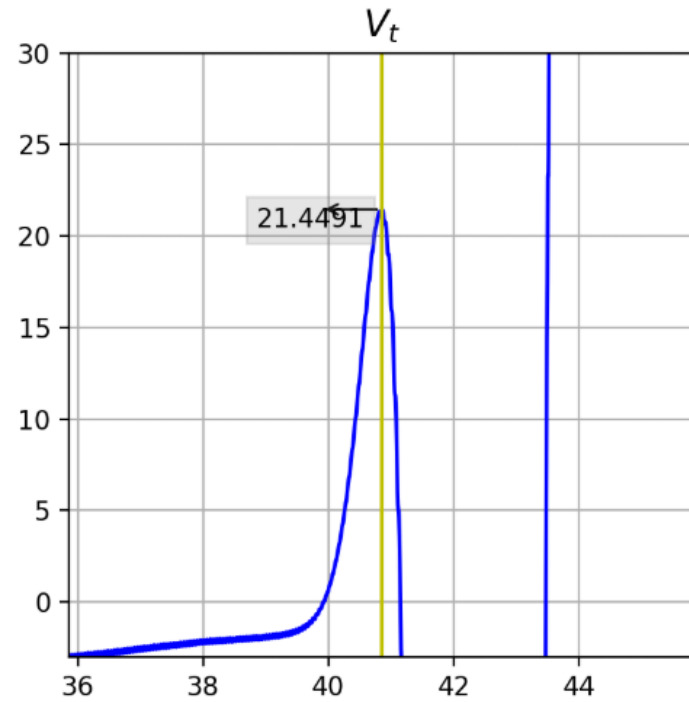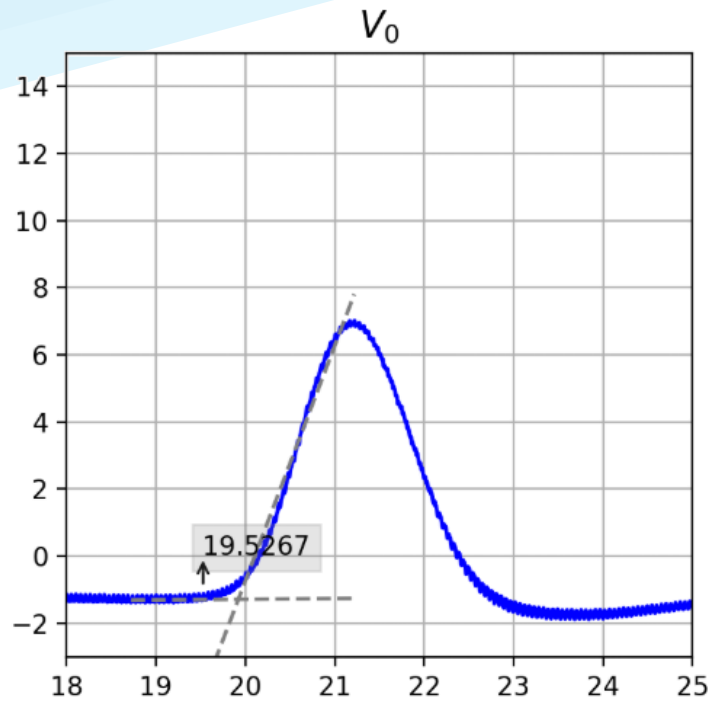$$odds\ ratio = \frac{DE/HE}{DN/HN} = \frac{20/380}{10/500}$$

# Quantile-quantile Plot (Q-Q Plot)

# Combination Chart

# Question Time

If you have any questions, please do not hesitate to ask me.

# The End

*Thank you for your attention ))*