

1.1 Data Attributes

(To be done over Week #2)

- (1) For the type of data below, state which would be considered continuous or discrete data types:

- (a) Number of patients in the wards *discrete*
- (b) The blood pressure of the patient *continuous*
- (c) The pulse of the patient *discrete*
- (d) The size of the patient's tumour *continuous*
- (e) The emergency room waiting time rounded to the nearest minute *discrete*

- (2) Which of the **NOIR** (Nominal, Ordinal, Interval, Ratio) scale would describe each of the data attributes below:

- (a) Rank in the Singapore Army *Ordinal*
- (b) Number of students registered for CZ4124 *Ratio*
- (c) Shoe sizes *Interval*
- (d) Jersey numbers in the Singapore national football team *Nominal*

	A	B	C	D	E	F	G
1	Student ID	Gender	Year of Birth	Attendance (weeks)	Score (%)	Grade	Completion Time (sec)
2	1	M	2002	12	85.5	A+	1600.31
3	2	F	2004	2	20	F	600.13
4	3	F	2003	13	55.5	C+	1800.00
5	4	M	1999	10	73.5	B+	1800.00
6	5	M	2000	11	65	B	1500.23
7	6	F	2004	8	70	B+	1700.00
8	7	F	2005	13	90.5	A+	900.34

Table 1 – Excel Table of Students Particulars and their performance in an online test

- (3) Which **NOIR** scale measure best represent each of the data category (column) in the Table 1?
- (4) Which measure(s) in the Table 1 (if any) can be considered a continuous data variable? Provide suitable justification for your answer.

Fifty percent of Android owners are under the age of 35

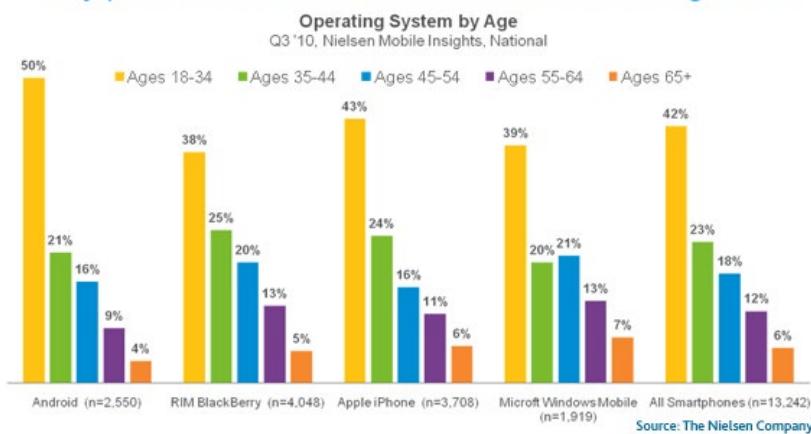


Figure 1 – Mobile OS by Age (From <https://www.nielsen.com/us/en/insights/article/2010/mobile-snapshot-smartphones-now-28-of-u-s-cellphone-market/>)

- (5) How many data dimensions have been simultaneously visualised in the clustered bar (column) chart shown in Figure 1? Draw a possible table (*partial table will do*) for this dataset and state which **NOIR** scale measure best represent each data dimension in your table.

1.2 Data Wrangling

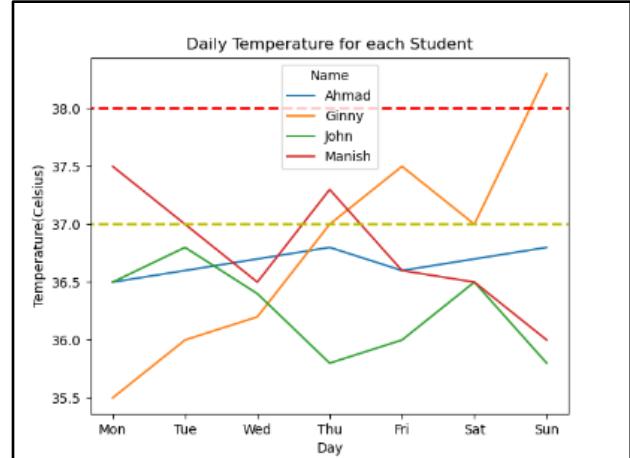
- (1) **Reshaping dataframes.** Figure 2(a) shows the table that can be extracted from the datafile “*Daily_Temperature.csv*” (available in NTULearn). Using the template file “*1-2-Daily_Temperature.py*”, plot the line chart shown in Figure 2(c). In order to get to the stage where you can directly plot the line chart shown, the dataframe created by the table shown in Figure 2(a) must be reshaped to that shown in Figure 2(b). Use the pandas **melt** and **pivot** functions to reshape the dataframe.

Name	Mon	Tue	Wed	Thu	Fri	Sat	Sun
John	36.5	36.8	36.4	35.8	36.0	36.5	35.8
Ahmad	36.5	36.6	36.7	36.8	36.6	36.7	36.8
Ginny	35.5	36.0	36.2	37.0	37.5	37.0	38.3
Manish	37.5	37.0	36.5	37.3	36.6	36.5	36.0

(a)

Name	Ahmad	Ginny	John	Manish
Day				
Mon	36.5	35.5	36.5	37.5
Tue	36.6	36.0	36.8	37.0
Wed	36.7	36.2	36.4	36.5
Thu	36.8	37.0	35.8	37.3
Fri	36.6	37.5	36.0	36.6
Sat	36.7	37.0	36.5	36.5
Sun	36.8	38.3	35.8	36.0

(b)



(c)

Figure 2 – (a) Daily temperature table from “*Daily_Temperature.csv*”. (b) The required dataframe format to plot the line chart shown in (c).

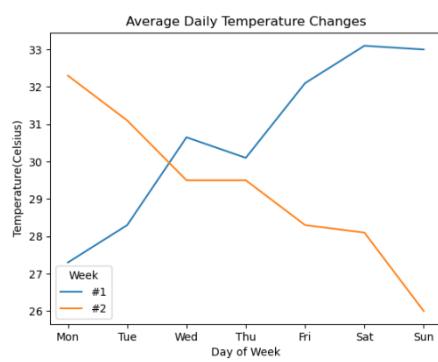
- (2) **Data wrangling exercise.** Figures 3(a) and 3(b) show the tables in the datafiles “*Temp_week_1.csv*” and “*Temp_week_2.csv*” respectively (available in NTULearn). Apply the appropriate data wrangling techniques (in pandas) to produce the two different line charts shown in Figures 3(c) and 3(d). Fill in the missing temperature data using the weekly mean temperature.

Challenge: If you can, fill in the missing data with a smoother interpolation from the two temporally nearest available data points.

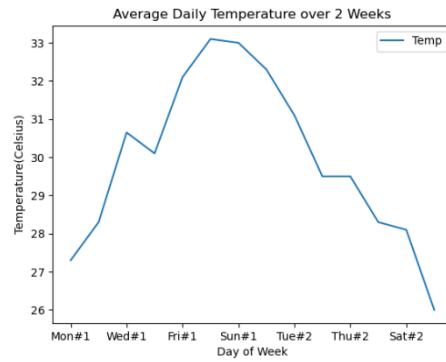
Day	Week	Temp
Mon	#1	27.3
Tue	#1	28.3
Wed	#1	
Thu	#1	30.1
Fri	#1	32.1
Sat	#1	33.1
Sun	#1	33.0

(b)

Day	Week	Temp
Mon	#2	32.3
Tue	#2	31.1
Wed	#2	29.5
Thu	#2	29.5
Fri	#2	28.3
Sat	#2	28.1
Sun	#2	26.0



(c)



(d)

Figure 3 – (a) Temperature tables for week #1 and (b) week #2. (c) Weekly temperature changes plotted as two different lines. (d) Single line plot of temperature changes over 2 weeks.

Appendix A – Listing of 1-2-Daily Temperature Python source file

```
# -*- coding: utf-8 -*-
"""
Created on Mon Aug  2 13:17:07 2021

CZ4124 Data Visualisation (Tutorial 1 Template)
@author: Put your name here

"""

import matplotlib.pyplot as plt
import pandas as pd

PlotWithPandas = True # you can plot either with Pandas or Matplotlib
#-----
# Read in Daily Temperature datafile into dataframe Temp
#-----
Temp = pd.read_csv('C:/Datasets/Daily_Temperature.csv')      # change to your
directory
print('\nTable read in\n',Temp,'\n')

#-----
# Use MELT to convert to long form and apply column names
#-----
# put in your code here

#-----
# Use PIVOT to convert to wide form with Names in each column
#-----
# put in your code here

#-----
# Use Pandas to plot line chart
#-----
if(PlotWithPandas):
    print('\nPlotting with Pandas')
    # put your Pandas plotting code here

#-----
# Use Matplotlib to plot line chart
#-----
else:
    print('\nPlotting with Matplotlib')
    # Alternatively, you can do the plot using Matplotlib or
    # any other Python plotting library
```


2.1 Visual Marks and Variables

(To be done over Week #3)

- (1) What are visual marks for items? List the three visual marks proposed by Jacque Bertin and describe the main differences in their basic property.
- (2) Name three perceptual channels (visual variables) proposed by Bertin and describe how they would alter the appearance of a line mark.
- (3) In the context of Bertin's visual marks and variables, describe what visual marks are employed in the visualisation shown in Figure 1 and what perceptual channels are employed to encode the different various information. Discuss the potential limitations of the perceptual channels employed.

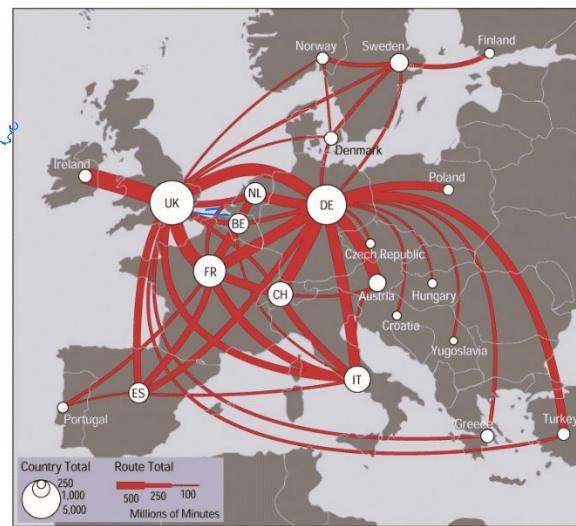


Figure 1 – Telecommunication Traffic Flow Map, © 2000 - TeleGeography, Inc.

Image taken from https://mappa.mundi.net/maps/maps_014/. For more details on the visualisation, check out the link.

2.2 Visual Encoding

- (1) Using your proposed visual encoding design, sketch a possible chart that will allow the effective visualisation of each of data tables show in Figures 2(a) to (c). The suggested visual mark for each chart has been indicated but you are free to choose an alternative. Consider carefully the attributes listed in each of the data table when designing your visual encoding and state any assumptions you made regarding the nature of these attributes that has influenced the design.
- (2) What potential visual quality problems may arise due to nature of the given data attributes? How do you proposed this problem can be addressed or reduced?

Lines

Living Members from Class of '45			
	Names	Gender	Age
1			
2			
3			
4			
5			
6			

(a)

Points

Orientation T-shirts - Comparing Halls					
No.	Gender	Weight	Height	Hall (A,B,C,D)	T-shirt Size
1					
2					
3					
:					
700					

(b)

Areas

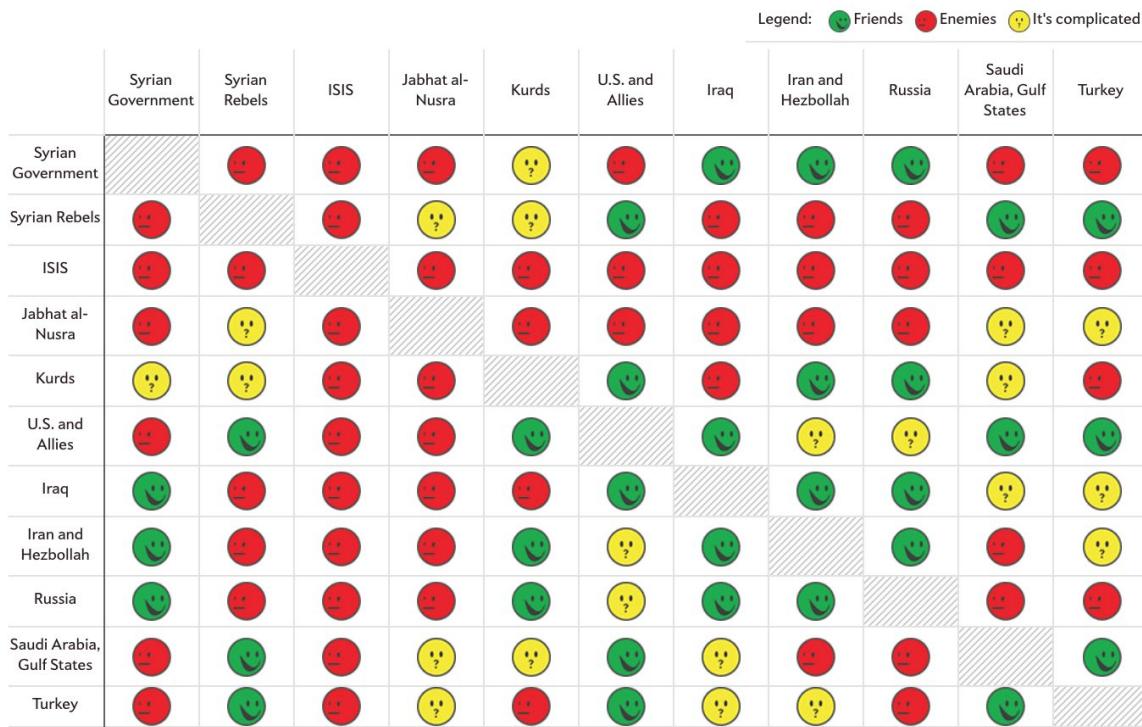
COVID-19 Infection in 2020			
Dates	Total Cases	Warded	ICU
1			
2			
3			
:			
365			

(c)

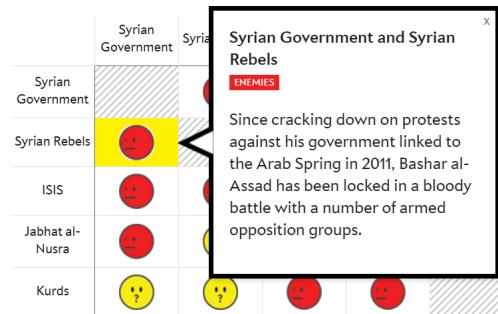
Figure 2 – Data table examples and the suggested visual mark that can use to design the chart.

2.3 Deconstruction Exercises

- (1) **Middle East Friendship Chart.** Figure 3(a) shows an interactive visualisation done by Joshua Keating and Chris Kirk (@Slate.com). It depicts the relationships between the active players involved in the civil war in Syria. If you click on any coloured emoji icon, a text description is displayed next to the icon describing the nature of the relationship between the two parties, as shown in Figure 3(b).



(a)



(b)

Figure 3 – Middle East Friendship Chart

(From http://www.slate.com/blogs/the_world/2014/07/17/the_middle_east_friendship_chart.html)

- Draw the data table for the relational data model used to create the visualisation shown in Figure 3. You can give appropriate names to the different attributes in your columns.
- What visual marks and variables (channel) are used to encode the different data depicted? *Color, Position, Shape*
- Describe if redundant coding has used in the visualisation and if so, discuss if it was done effectively.
- Describe if the choice of colours used in visual encoding is appropriate.

2.4 Critical Eye

(1) **BBC News – Fortunes of Chinese Family.** On 15 Oct 2017, BBC News published an article entitled **Five charts about the fortunes of the Chinese family**. In this article, 5 charts were presented on various statistical facts about the state of the Chinese family in China. These charts are shown in Figure 4.

Look at these charts and the visual encoding used in each case. Critic these visuals (infographics). Comment if they were done well (i.e. the information intended to be conveyed was effectively communicated with the choice of visual encoding). Suggest if any of these charts could be improved visually to allow it to express the intended message more clearly and quickly.

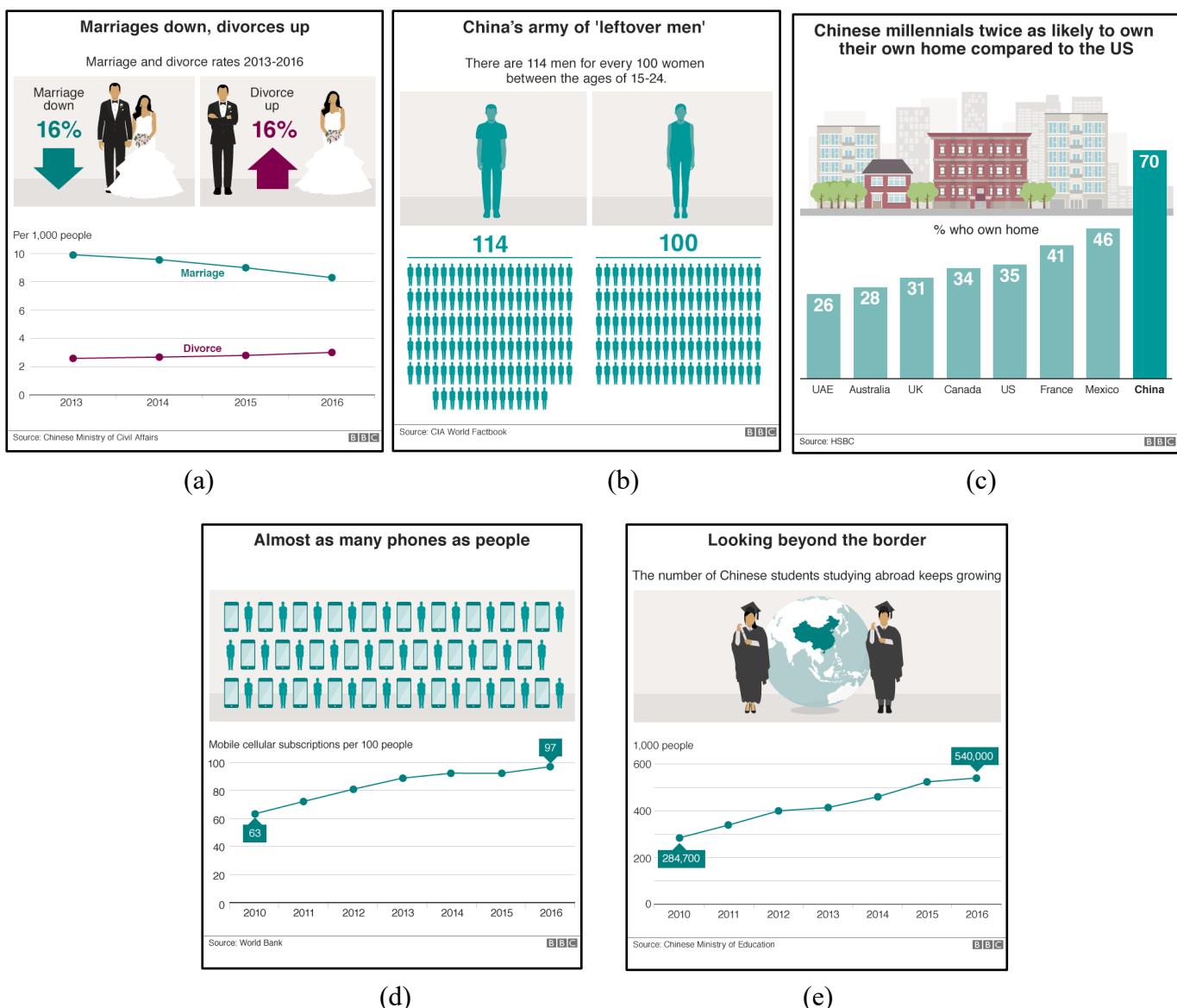


Figure 4 – Five charts from a BBC News article (Five charts about the fortunes of the Chinese family)
(Image from <https://www.bbc.com/news/world-asia-china-41424041>)

Optional Challenge (Give this a try)

(2) **Euro Cup 2008.** Duch, Waitzman and Amaral analysed three knockout-phase matches of the Spanish team using the vast amount of statistical information that was published online by UEFA. Using techniques from social network analysis, they produced the visualisation shown in Figure 5. Based on further descriptions given in their PLOS ONE open-access paper accessible at <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0010937>, answer the following questions:

- Identify as many types of data variables as possible that have been visually represented in the visualisation shown in Figure 4. For each, state its scale of measure (N, O or Q).
- Describe what perceptual channel (visual variable) has been employed to encode the values of each of these data variables.
- Describe if the choice of colours used in visual encoding is appropriate. If not, suggest how it could be improved.

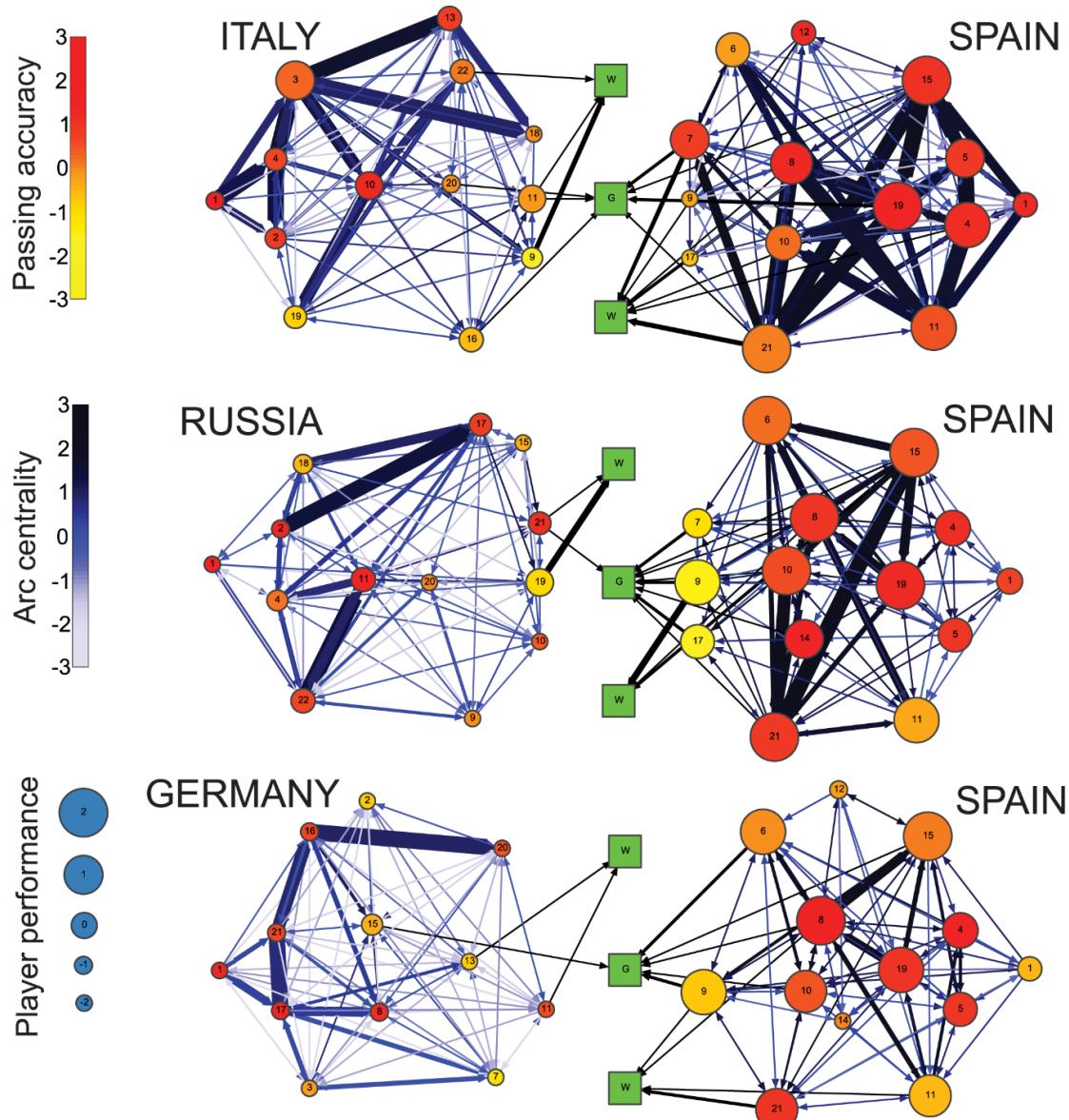


Figure 5 – Visualisation of the three knockout-phase matches of the Spanish team.
(Image from <https://journals.plos.org/plosone/article/figure?id=10.1371/journal.pone.0010937.g005>)

2.1 TUT 2

- points
- lines
- areas

z - position
for line

- size ✓
- value
- texture
- colour
- orientation
- shape

3

points

lines

area

position

size ✓

texture

colour

2.2 1)

Lines

Living Members from Class of '45

	Names	Gender	Age
1			
2			
3			
4			
5			
6			

(a)



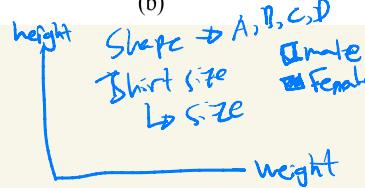
◻ male
● female

Points

Orientation T-shirts - Comparing Halls

No.	Gender	Weight	Height	Hall (A,B,C,D)	T-shirt Size
1					
2					
3					
:					
700					

(b)

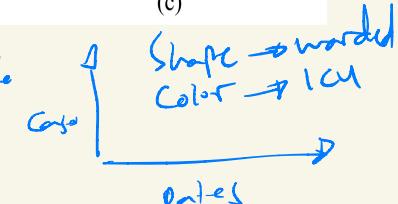


Areas

COVID-19 Infection in 2020

Dates	Total Cases	Warded	ICU
1			
2			
3			
:			
365			

(c)



3.1 Line and Bar Charts

(To be done over Weeks #4 & #5)

- (1) Table 1 shows the transaction volume received and processed over the last 12 months of Jan to Dec. This data table can be found in the file “*hire2FTE.csv*”. Using the Matplotlib and/or Seaborn Python visualisation library, create the cluster bar chart shown in Figure 1(a). Do note that you may have to perform some data wrangling (as in Tutorial #1) to prepare your dataframe for plotting the bar chart shown.

You are to annotate the pertinent information shown, which include the following:

- a) Line marker showing two workers quit in the month of May.
- b) The actual transaction volumes received and processed for each month.
- c) The title of the plot is “**Please approve our request to hire 2 new workers**”
- d) The *x* and *y* axis labels.

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Transactions	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
2	Received	160	170	240	140	175	155	130	200	160	140	150	175
3	Processed	160	170	240	140	175	150	125	170	135	130	125	150

Table 1 – The data table in the csv data file “*hire2FTE*”

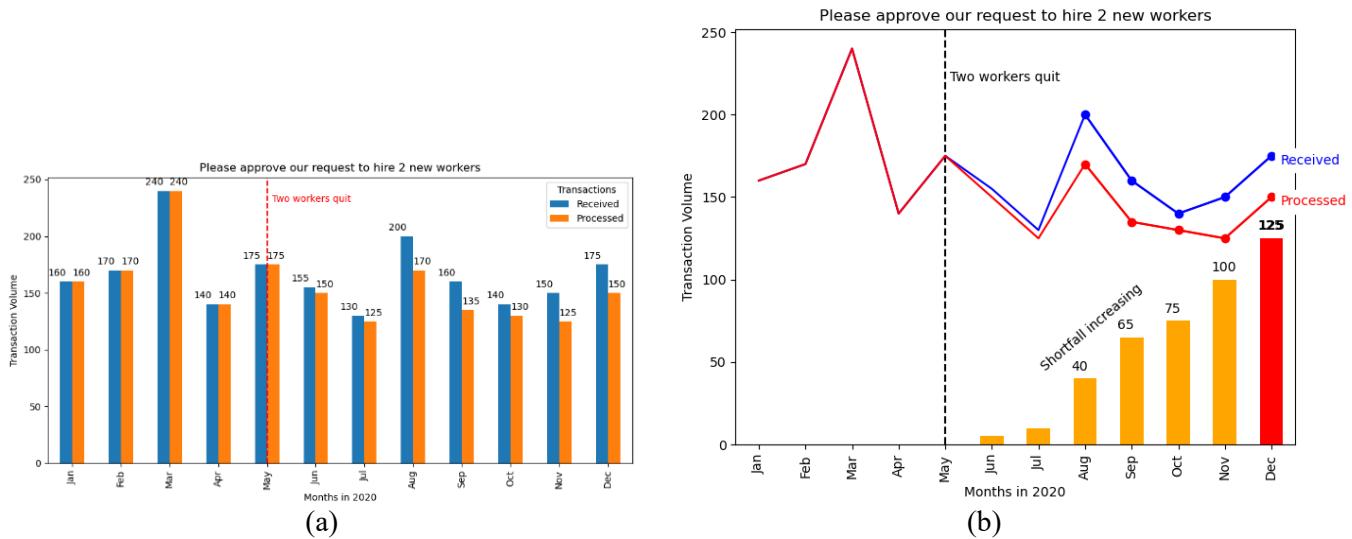


Figure 1 – (a) Cluster bar chart showing the monthly volumes of transaction receive and processed. (b) An equivalent combination of line and bar plots showing the same information.

- (2) The purpose of your data visualisation is to present a case to your upper management that you need to quickly hire two additional full-time employees (FTE) as your team’s ability to keep up with transaction volume received is falling way behind. Comment on the main reason why the bar chart in Figure 1(a) is not suitable for this purpose.
- (3) Create the two red and blue line plots shown in Figure 1(b), along with the various marker and text annotations shown. Comment on their effectiveness, including choice of colours, marker placement, text annotation, highlights, etc. Suggest ways to further improve (if any).
- (4) **Challenge (Optional):** Superimpose the increasing shortfall bar plots over your two line plots. The shortfall is computed by cumulating the monthly differences between the transactions received and processed. Include the annotations and highlights as shown in Figure 1(b) and comment on their effectiveness.

3.2 The Perils of Line Smoothing

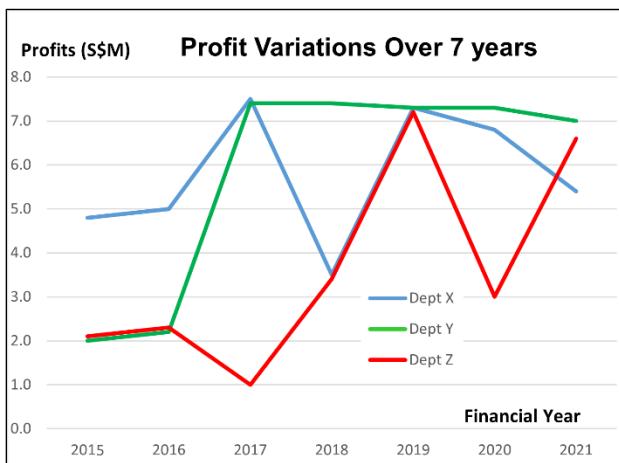
(1) Smoothing data series is a method to remove local variations (noise) in data so that a more general trend can be observed in the data variations. However, there are situations when such smoothing operations actual distort the truth about what the data is actually telling us. First, study the profit figures in Table 2 and answer the following questions:

- Which department achieved the highest ever annual profit over the last 7 years and in which year was this? *X (7.5 in 2017)*
- Did Dept Z ever had higher annual profits compared to Dept X? If so, over which period was this? *Yes (2021)*
- Describe the difference in profits differences between Dept Y and Z during the earlier days of this company (i.e. years 2015 and 2016). *0.1*
- Year 2017 to 2019 were years of change. Did any department managed to maintain their profits? *Only dept Y keeping steady between 7.3 to 7.4, no change at all in 2017 to 2019*
- Describe the relative performance of the three departments in year 2019.

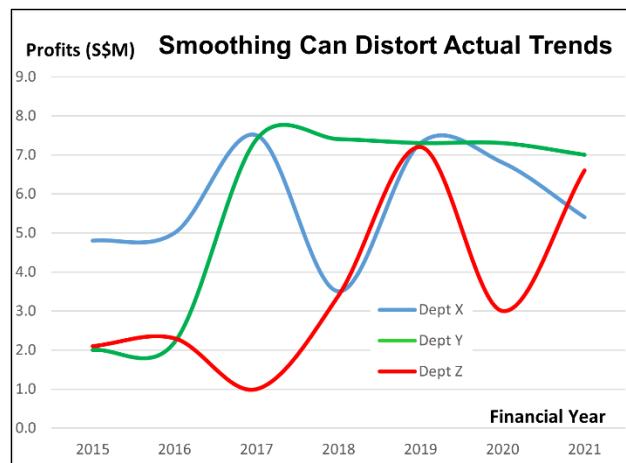
Profits in last 7 years (\$S M)			
Year	Dept X	Dept Y	Dept Z
2015	4.8	2.0	2.1
2016	5.0	2.2	2.3
2017	7.5	7.4	1.0
2018	3.5	7.4	3.4
2019	7.3	7.3	7.2
2020	6.8	7.3	3.0
2021	5.4	7.0	6.6

Table 2 – Department Profits over 7 years.

(2) Second, study the two line charts in Figure 2 carefully and described all the misleading trends created due to the application of the Excel smoothed line function in Figure 2(b)?



(a)



(b)

Figure 2 – Excel line chart plots of the departmental profit variations in Table 1. (a) The three line plots of the actual profits of each departments X, Y and Z. (b) The line chart plotted using the “Smoothed line” option in Excel Line Chart.

(3) Challenge (Optional):

Using the Excel file “dept profits.xlsx” provided, create a clustered bar chart shown in Figure 2(d) to help you compare the annual profits between the three departments over the period of 2015 to 2021. You can do this quickly in Excel.

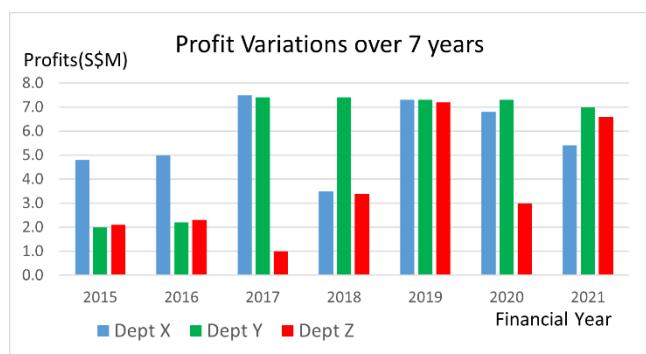


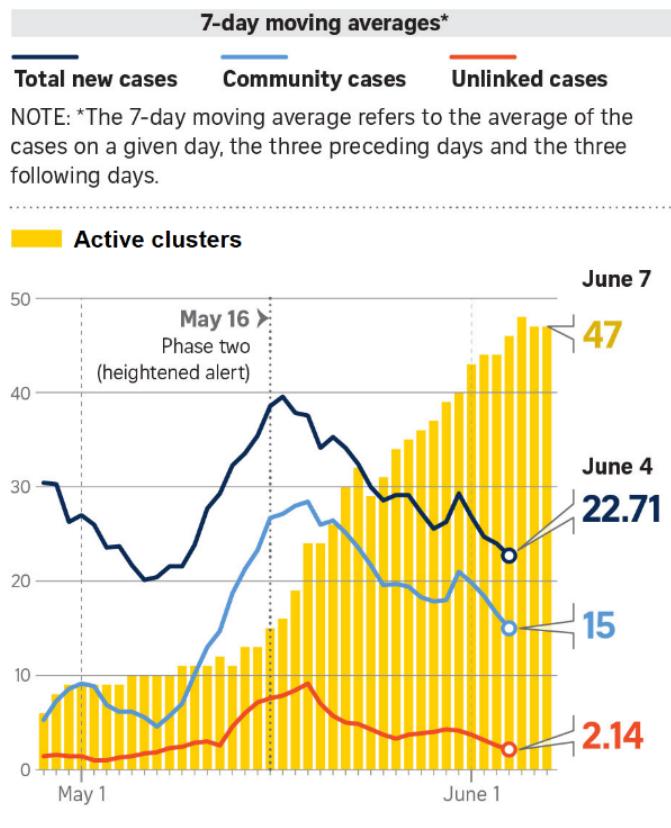
Figure 2(d)

3.3 Critical Eye

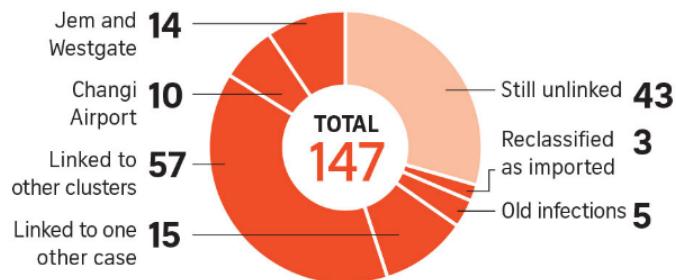
(1) **Communicating Singapore's COVID-19 Situation.** Figure 3 shows an infographic taken from The Straits Times. It describes a snapshot of the COVID-19 infection situation in Singapore from the start of May till the day of publication on 8 June 2021.

- Why do you think both line and bar charts were used in visualising the virus infection data shown in the upper portion of Figure 3?
- Three lines of different colours were used to depict the total new cases, community cases and unlinked cases. Would a stacked area chart be more appropriate since these data values by their very nature, will never overlap? Give reasons for your answer.
- What about the bar chart in the background? Would it be better to replace it with an area chart of the same colour?
- Why was a marker line drawn over the line charts at May 16? What can you say about the perceptual channel used on this marker line?
- What story do you think is being communicated with the inclusion of this marker at May 16?
- Why was a donut chart used to visualize the unlinked cases?
- Could the donut chart be replaced with a pie chart? If so, what changes would you make to the visual design.
- Should each slice of the donut chart be encoded with a different colour to improve visual clarity? Is there any particular reason why only two colours have been used in the donut chart?

Singapore's virus situation since April 28



Cases initially reported as unlinked



Source: MINISTRY OF HEALTH
STRAITS TIMES GRAPHICS

Figure 3 – COVID-19 Infection Situation in Singapore

(Image modified from The Straits Times on 8 June 2021 at <https://www.straitstimes.com/singapore/health/unlinked-covid-19-cases-connection-to-clusters-found-within-days-in-singapore>)

TUT 3

3.1 Line and Bar Charts

(1) Plotting Line Chart

Plot the two lines & dot markers

```
import pandas as pd
import matplotlib.pyplot as plt
```

```
# Read in csv file and separate Rec and Proc
Transaction = pd.read_csv('C:/Datasets/HrCapTE.csv')
Receive = Transaction.loc[0, 'Jan': 'Dec']
Process = Transaction.loc[1, 'Jan': 'Dec']
MontOff = Receive['Proc']
ShortFall = Process['Rec']
```

```
# plot two lines for receive and Process
plt.plot(Receive, color='blue') # plot receive line in blue
plt.plot(Process, color='red') # plot process line in red
```

draw markers

```
plt.plot(Receive[Aug : Dec], marker='o', color='blue') # blue dots
plt.plot(Process[Aug : Dec], marker='o', color='red') # red dots
```

3.1 Line and Bar Charts

(1) Plotting Line Chart

Compute and plot shortfall-bars

```
import pandas as pd
import matplotlib.pyplot as plt
```

```
# Read in csv file and separate Rec and Proc
Transaction = pd.read_csv('C:/Datasets/HrCapTE.csv')
Receive = Transaction.loc[0, 'Jan': 'Dec']
Process = Transaction.loc[1, 'Jan': 'Dec']
MontOff = Receive['Proc']
ShortFall = Process['Rec']
```

Calculate the monthly difference into shortfall

```
Total = 0
for i in range(12):
    Total += MontOff[i]
    Total -= Process[i]
    Total -= Receive[i]
    Total -= ShortFall[i]
    print(Total)
    Total = 0
```

Draw shortfall bar plot using Pandas

```
fig1 = ShortFall.plot.bar(ylabel="Shortfall")
```

3.1 Line and Bar Charts

(1) Plotting Line Chart

Highlight last shortfall bar in red

```
import pandas as pd
import matplotlib.pyplot as plt
```

```
# Read in csv file and separate Rec and Proc
Transaction = pd.read_csv('C:/Datasets/HrCapTE.csv')
Receive = Transaction.loc[0, 'Jan': 'Dec']
Process = Transaction.loc[1, 'Jan': 'Dec']
MontOff = Receive['Proc']
ShortFall = Process['Rec']
```

highlight the last shortfall bar

```
LastshortFall = ShortFall.copy()
LastshortFall['Dec'] = 0
fig1 = LastshortFall.plot.bar(color='red') # plot all red
```

```
fig1 = LastshortFall.plot.bar(color='orange') # turn Jan-Nov back to orange
```

3.1 Line and Bar Charts

(1) Plotting Line Chart

Annotate values of shortfall bars

```
import pandas as pd
import matplotlib.pyplot as plt
```

```
# Read in csv file and separate Rec and Proc
Transaction = pd.read_csv('C:/Datasets/HrCapTE.csv')
Receive = Transaction.loc[0, 'Jan': 'Dec']
Process = Transaction.loc[1, 'Jan': 'Dec']
MontOff = Receive['Proc']
ShortFall = Process['Rec']
```

Label shortfall bar numbers

```
LabelCount = 5
LabelMonth = 'Aug'
Index = 0
LabelMonth = 'Aug' + ' ' + 'Dec'
Index = LabelMonth
textXPos = Index
textYPos = Index + 1
plt.text(textXPos, textYPos, Index, color='black')
plt.text(textXPos, textYPos, Index, color='black', fontweight='bold')
```

3.1 Line and Bar Charts

(1) Plotting Line Chart

Label title, axes and shortfall bars

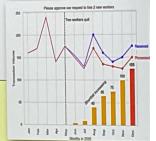
```
import pandas as pd
import matplotlib.pyplot as plt
# Read in csv file and separate Rec and Proc
Transaction = pd.read_csv('C:/Datasets/HrCapTE.csv')
Receive = Transaction.loc[0, 'Jan': 'Dec']
Process = Transaction.loc[1, 'Jan': 'Dec']
MontOff = Receive['Proc']
ShortFall = Process['Rec']

# plot two lines for receive and Process
plt.plot(Receive, color='blue') # plot receive line in blue
plt.plot(Process, color='red') # plot process line in red

# draw markers
plt.plot(Receive[Aug : Dec], marker='o', color='blue') # blue dots
plt.plot(Process[Aug : Dec], marker='o', color='red') # red dots

# rotate shortfall annotated text
plt.text(45, ShortFall[Increasing], color='black', rotation=45)

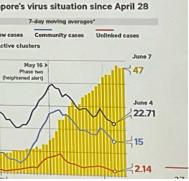
# add axes labels and title
plt.title('Please approve our request to hire 2 new workers') # plot title
plt.xlabel('Months in 2020') # x-axis label
plt.ylabel('Transaction Volume') # y-axis label
```



3.3 Critical Eye

(1) Communicating SG's COVID Situation

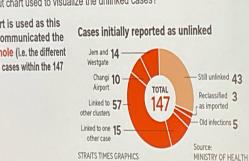
- Why do you think both line and bar charts were used in visualising the virus infection data?
- The bar and line describe different types of data.
- The bar describes the number of active clusters.
- The line describes number of individual cases.
- Using both helps make this distinction clear and avoid confusion.



3.3 Critical Eye

(1) Communicating SG's COVID Situation

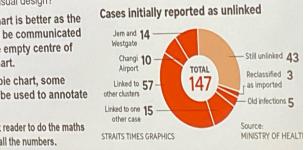
- What was a donut chart used to visualize the unlinked cases?



3.3 Critical Eye

(1) Communicating SG's COVID Situation

- Could the donut chart be replaced with a pie chart? If so, what changes would you make to the visual design?



3.3 Critical Eye

(1) Communicating SG's COVID Situation

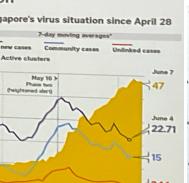
- What about the bar chart in the background? Would it be better to replace it with an area chart of the same colour?



3.3 Critical Eye

(1) Communicating SG's COVID Situation

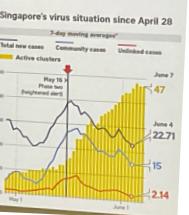
- Possible but not as effective as bar chart.
- Individual bars capture exact value for each day and the abrupt changes in data values much better.
- Texture of bar chart gives good contrast to the 3 line plots.



3.3 Critical Eye

(1) Communicating SG's COVID Situation

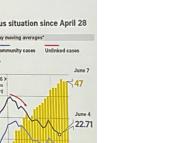
- Why was marker line drawn at May 16?
- What can you say about the perceptual channel used on this marker line?
- The line at May 16 marks an important event (start of the imposition of Phase 2 - Heightened alert).
- This line is dotted (texture perceptual channel) instead of solid so it is visually distinct from the three lines use for the line chart.



3.3 Critical Eye

(1) Communicating SG's COVID Situation

- What story do you think is being communicated with the inclusion of this marker at May 16?
- By drawing this time marker over the line charts, it communicates the effectiveness of the heightened alert.
- Line charts show clear changing trend of falling infection cases soon after the line at May 16.

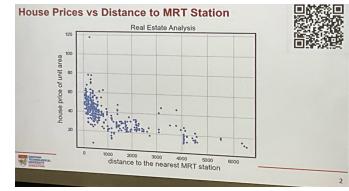


more effective to use darker colour to highlight key segment

3.4 Analysing Relationships

(1) **What Affects Property Prices?** The data file “*Real estate.csv*” provides data on 414 property transactions. The transaction price for each property sold is quoted based on cost per unit area, which means it gives the general value of the property and is not affected by its size. Associated with each transaction is also listed the following characteristics of the property that is of concern to us, namely:

- a) Unit area price of the property (“**house price of unit area**”)
- b) Age of the property (“**house age**”)
- c) Distance from a MRT station (“**distance to the nearest MRT station**”)
- d) Number of convenient stores nearby (“**number of convenience stores**”)



Using an appropriate chart and relationship model fitting technique (e.g. linear regression, polynomial regression, lowess, etc) visualise the general relationships between the various variables list above.

Real estate dataset from: <https://www.kaggle.com/quantbruce/real-estate-price-prediction?select=Real+estate.csv>

(2) **Distance to MRT** - What can you say about the general price of a property with respect to its distance from the nearest MRT station? What do you think is the reason for this relationship? Is this relationship a simple linear one or is it more complex? Tell this story with evidence from your visualizations.

Food for thought: <https://www.propertyguru.com.sg/property-guides/mrt-effect-on-property-prices-39498> (accessed in June'21)

(3) **Age of the property** - What can you say about the general price of a property with respect to its age? Does old really mean cheap? What is the reason for the relationship (Hint: Does the issue in part (2) has something to do with this? Make your case with evidence from your visualizations).

(4) **Number of convenience stores** – Do note that the data on **number of convenience stores** near the property is discrete and ordinal, as such it should be treated differently from the other data categories available in the data set (Hint: Plot the number of convenience stores on the x-axis so you can do multiple horizontal distribution visualisations using appropriate Seaborn plots).

- (a) What can you say about the general price of a property with respect to the number of convenience stores in the neighbourhood? Make your case with evidence from your visualizations.
- (b) What kinds of neighbourhood has many convenience stores? Do you think you can use the data in *Real estate* to do this analysis? If you can, in the light of this new analysis, what can you say about the case you made earlier in part 4a?

(5) **Faceting** - Relationships influencing a primary variable like the house price is multi-faceted and making hasty conclusions with a single visualisation should be avoided without first visualising and analysing the relationships between the other variables.

- (a) Explore Seaborn’s many powerful methods to create grids of multiple plots to explore the various relationships (checkout: https://seaborn.pydata.org/tutorial/axis_grids.html).
- (b) What do you think is the primary factor influencing the general price of a property based on the data available in *Real estate.csv*? Provide the visualisations to support your view.

3.5 Analysing Distributions

(1) **Factors Affecting Students Academic Performance.** The data file “*Students performance.csv*” provides data on the math, reading and writing scores of 1000 students from five different ethnic groups A to E. Other information like gender, parent’s education and whether they have completed their test preparation course are also given. Analyse these data with appropriate distribution plots. In some cases, you will need to **wrangle the data** into suitable format before you can use the various Seaborn distribution plots to visualise and compare the distributions.

Students performance dataset from: <https://www.kaggle.com/spscientist/students-performance-in-exams>

(2) **Gender Differences** - What can you say about the relative average performance differences between male and female students across the three different subjects? What subjects do girls generally do better in and what subjects do the boys do better in? Show this analysis using appropriate comparative distribution plots (*Hint:* consider using multiple violin plot, with its left/right split feature).

(3) **Toughest Subject** – Which of the three subjects did the entire cohort performed the worst in? Show this analysis using appropriate comparative distribution plots and annotate the median value for each of the subjects into your plot (*Hint:* consider using multiple box plots).

(4) **Test Preparation Course** – Does completing the test preparation course help students to do better in the three different subjects? Which subject shows the most improvement when the preparation course is completed? Show this analysis using appropriate comparative distribution plots and annotate the median values for both completion and non-completion of preparation course for each of the subjects into your plot to make it easy to visualize the different medians.

(5) **Performance across Ethnic Groups** – How did the different ethnic groups A to E performed in their tests? Which group did the best overall in all three subjects? (*Hint:* consider using a series of overlapping kernel density estimate plots).

(6) **Subject Performance across Ethnic Groups** – Were there performances differences in the three different subjects across the five different groups A to E (e.g. did one group do the best in math and another group did the best in reading, or did one group did the best in all three subjects)? (*Hint:* consider using a series of coloured box plots).

(7) **Influence of Parent’s Education on Students’ Performance** – Did the parent’s educational background have an influence on their children’s performance in general (i.e. all their subjects combined)? If they did, which of these educational backgrounds seems to result in the best performance and which resulted in the worst?

(8) **Influence of Parent’s Education on Gender** – If the parent’s educational background have an influence of the students’ performance, was this influence more pronounce for the male or female students?

Optional Challenge

(Additional optional dataset to explore)

3.6 Analysing the Data

- (1) **What Happened to the Pandemic in Malaysia?** A data file entitled “owid MYS.csv” has been provided for you, which is a subset of COVID data related to the COVID pandemic cases in Malaysia extracted on 9 June 2021 from the Our World in Data website at <https://github.com/owid/covid-19-data/tree/master/public/data>. Study the data table and locate the column labelled “**new_cases**”.
- (2) Using the Seaborn and any other relevant Python libraries, wrangle the data table given to you and plot the bar chart for the **new_cases** over the 500 days starting on the date ‘**25/1/2020**’.
- (3) Superimpose over the bar chart an appropriate smoothed line so that you can visualize the general changing trends of number of new confirmed COVID-19 cases over the 500 days.
- (4) Based on some of the relevant COVID-19 related news and web links (and any others you think relevant), annotate on your plot (i.e. using appropriate visual markers and text) any relevant events that you think may have influenced the changing infection trends in your chart.

References to some COVID-19 relevant events in Malaysia:

General info - https://en.wikipedia.org/wiki/COVID-19_pandemic_in_Malaysia

27 Feb 2020 - https://en.wikipedia.org/wiki/2020_Talibqi_Jamaat_COVID-19_hotspot_in_Malaysia

26 Sep 2020 - <https://www.straitstimes.com/asia/se-asia/malaysias-pm-muhyiddin-admits-sabah-state-polls-in-sept-caused-current-covid-19-wave>

23 Dec 2020 - <https://www.straitstimes.com/asia/se-asia/malaysia-has-identified-new-covid-19-strain-similar-to-one-found-in-3-other-nations>

16 Jan 2021 - <https://asia.nikkei.com/Spotlight/Coronavirus/Malaysia-s-Top-Glove-reports-COVID-19-outbreak-at-four-factories>

13 Apr 2021 - <https://www.straitstimes.com/asia/se-asia/malaysia-relaxes-some-coronavirus-curbs-as-fasting-month-arrives>

2 May 2021 - <https://www.bangkokpost.com/world/2109227/malaysia-reports-first-case-of-indian-covid-19-variant>

Note: Some of the web links may be out of date and no longer accessible. You can infer the related event from the title of the web link.

3.4 Analysing Relationships

(2) Distance to MRT station

Q13. MRT - Which factors are substantiated by the scatter plots?

- The closer the house is to an MRT station, the higher its price. (67% correct)
- The closer the house is to an MRT station, the lower its price. (29% correct)
- The closer the house is to an MRT station, the more expensive it is. (17% correct)
- The closer the house is to an MRT station, the less expensive it is. (17% correct)
- The closer the house is to an MRT station, the more it costs. (17% correct)
- The closer the house is to an MRT station, the less it costs. (17% correct)

Legend: Not hot or warm or empty
House price > 1M
House price < 1M
House price > 2M
House price < 2M
House price > 3M
House price < 3M
House price > 4M
House price < 4M
House price > 5M
House price < 5M
House price > 6M
House price < 6M
House price > 7M
House price < 7M
House price > 8M
House price < 8M
House price > 9M
House price < 9M
House price > 10M
House price < 10M
House price > 11M
House price < 11M
House price > 12M
House price < 12M
House price > 13M
House price < 13M
House price > 14M
House price < 14M
House price > 15M
House price < 15M
House price > 16M
House price < 16M
House price > 17M
House price < 17M
House price > 18M
House price < 18M
House price > 19M
House price < 19M
House price > 20M
House price < 20M
House price > 21M
House price < 21M
House price > 22M
House price < 22M
House price > 23M
House price < 23M
House price > 24M
House price < 24M
House price > 25M
House price < 25M
House price > 26M
House price < 26M
House price > 27M
House price < 27M
House price > 28M
House price < 28M
House price > 29M
House price < 29M
House price > 30M
House price < 30M
House price > 31M
House price < 31M
House price > 32M
House price < 32M
House price > 33M
House price < 33M
House price > 34M
House price < 34M
House price > 35M
House price < 35M
House price > 36M
House price < 36M
House price > 37M
House price < 37M
House price > 38M
House price < 38M
House price > 39M
House price < 39M
House price > 40M
House price < 40M
House price > 41M
House price < 41M
House price > 42M
House price < 42M
House price > 43M
House price < 43M
House price > 44M
House price < 44M
House price > 45M
House price < 45M
House price > 46M
House price < 46M
House price > 47M
House price < 47M
House price > 48M
House price < 48M
House price > 49M
House price < 49M
House price > 50M
House price < 50M
House price > 51M
House price < 51M
House price > 52M
House price < 52M
House price > 53M
House price < 53M
House price > 54M
House price < 54M
House price > 55M
House price < 55M
House price > 56M
House price < 56M
House price > 57M
House price < 57M
House price > 58M
House price < 58M
House price > 59M
House price < 59M
House price > 60M
House price < 60M
House price > 61M
House price < 61M
House price > 62M
House price < 62M
House price > 63M
House price < 63M
House price > 64M
House price < 64M
House price > 65M
House price < 65M
House price > 66M
House price < 66M
House price > 67M
House price < 67M
House price > 68M
House price < 68M
House price > 69M
House price < 69M
House price > 70M
House price < 70M
House price > 71M
House price < 71M
House price > 72M
House price < 72M
House price > 73M
House price < 73M
House price > 74M
House price < 74M
House price > 75M
House price < 75M
House price > 76M
House price < 76M
House price > 77M
House price < 77M
House price > 78M
House price < 78M
House price > 79M
House price < 79M
House price > 80M
House price < 80M
House price > 81M
House price < 81M
House price > 82M
House price < 82M
House price > 83M
House price < 83M
House price > 84M
House price < 84M
House price > 85M
House price < 85M
House price > 86M
House price < 86M
House price > 87M
House price < 87M
House price > 88M
House price < 88M
House price > 89M
House price < 89M
House price > 90M
House price < 90M
House price > 91M
House price < 91M
House price > 92M
House price < 92M
House price > 93M
House price < 93M
House price > 94M
House price < 94M
House price > 95M
House price < 95M
House price > 96M
House price < 96M
House price > 97M
House price < 97M
House price > 98M
House price < 98M
House price > 99M
House price < 99M
House price > 100M
House price < 100M

3.4 Analysing Relationships

(3) Age of Property

Q14. Age - Which factors are substantiated by these scatter plots?

- Property price increases linearly with age of the property. (37% correct)
- Age of the property is negatively distributed below 0 to 45 years. (17% correct)
- There is no strong correlation between property price and age. (88% correct)
- There is a strong positive correlation between property price and age. (67% correct)
- There is a strong negative correlation between property price and age. (50% correct)

3.4 Analysing Relationships

(2) Distance to MRT station

Property price and distance to MRT?

- Linear or more complex? Give reasons.

The **regplot** linear regression model (order=1) shows a clear fall in property price when further away from MRT station. Why?

- The accessibility to convenient transportation could explain this trend of why people willing to pay more to be closer to a MRT station.

regplot (order=1)

lowess (frac=0.4)

replot (order=1)

distance to the nearest MRT station

house price of unit area

3.4 Analysing Relationships

(2) Distance to MRT station

import seaborn as sns

```
Xtrain = MRT
Ytrain = Price
PolyOrder = 1
Fig1 = sns.regressionplot(data=Property, x=Xtrain, y=Ytrain,
                          order=PolyOrder,
                          line_kws={'color': 'red'},
                          scatter_kws={'s': 20, 'color': 'cyan'})
```

of why people willing to pay more to be closer to a MRT station.

regplot (order=1)

lowess (frac=0.4)

replot (order=1)

distance to the nearest MRT station

house price of unit area

3.4 Analysing Relationships

(2) Distance to MRT station

import seaborn as sns

```
from statsmodels.nonparametric.smoothers_lowess import lowess
Xtrain = MRT
Ytrain = Price
PolyOrder = 1
Fig1 = sns.regplot(data=Property, x=Xtrain, y=Ytrain,
                   order=PolyOrder,
                   line_kws={'color': 'red'},
                   scatter_kws={'s': 20, 'color': 'cyan'})
```

Fractal = 0.4

XVal = Property[Xtrain] # Get X values
YVal = Property[Ytrain] # Get Y values
LowessVal = lowess(XVal, YVal, frac=Fractal)
Fig1 = sns.lineplot(x=LowessVal[:,0], y=LowessVal[:,1],
 linewidth=2.0, color='green')

regplot (order=1)

lowess (frac=0.4)

distance to the nearest MRT station

house price of unit area

3.4 Analysing Relationships

(3) Age of Property

Property price and property age?

- Old = cheap? Give reasons.

The **lowess** (frac=0.5) shows that property price drops with increasing age of the property but as the age of the property crosses about 20 years, its value picks up again. Why?

- But where are these really old properties?

lowess (frac=0.5)

Real Estate Analysis

house price of unit area

house age

3.4 Analysing Relationships

(3) Age of Property

Property price and property age?

- Old = cheap? Give reasons.

Let's examine relationship between the property age and its proximity to the MRT.

- Many older homes past the 20-year mark are mostly near MRT stations. We saw earlier that this increases property price.

lowess (frac=0.5)

Real Estate Analysis

house price of unit area

house age

3.4 Analysing Relationships

(4) Number of convenience stores (ordinal values)

(a) Property price and no. of stores nearby?

- Violin plots used together with lowess (frac=0.5) to see the relationships between number of stores and property prices.
- The violin plots seem to suggest that the property prices generally go up with the increasing number of convenience stores.
- But the steepest rate of increase is between 3 and 5 stores nearby.

Violin point

wide spread (low spread)

wide spread (high spread)

house price of unit area

house age

3.4 Analysing Relationships

(4) Number of convenience stores (ordinal values)

(b) What kinds of neighbourhood has many convenience stores?

- Plot the Distance to MRT vs no. of convenience stores.
- The trend reveals that the large no. of stores are mainly around the MRT stations.
- This could imply that the property price is mainly affected by its proximity to the MRT stations instead of just the no. of stores nearby.

Distance to MRT

house price of unit area

number of convenience stores

3.4 Analysing Relationships

(5) Faceting

(a) Use multi-faceted plot to compare multiple factors simultaneously.

- The Seaborn **PairGrid** function can be used.
- Then apply **regplot** to the off-diag plots, with order=2.
- And **histplot** to the diag plots.

pair grid

regression plot

dot attributes

line attributes

polynomial order = 2

no confidence interval

add histogram (diagonal)

orange bars

Pair Grid plot

y-axis attributes

x-axis attributes

scale display size

regression plot

line attributes

dot attributes

polynomial order = 2

no confidence interval

add histogram (diagonal)

orange bars

Fig1 = sns.PairGrid(Property,
 y_vars=[Price, Store, MRT, Age],
 x_vars=[Price, Store, MRT, Age],
 height=2)

Fig1.map_offdiag(msn.regressionplot,
 line_kws={'color': 'red'},
 scatter_kwds={'s': 20},
 order=2,
 ci=None)

Fig1.map_diag(msn.histplot,
 color='orange')

3.4 Analysing Relationships

(5) Faceting

(b) What is the primary factor influencing the general price of property?

- Main factor seems to be proximity to MRT station.
- Influence on property prices of store age and no. of stores have a corresponding opposite trend when plotted against distance to MRT station.

MRT

Store

Age

pair grid

regression plot

dot attributes

line attributes

polynomial order = 2

no confidence interval

add histogram (diagonal)

orange bars

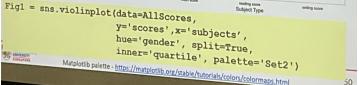
3.5 Analysing Distributions

(2) Gender Differences

Relative performance between male and females?

Solution:

- Use split violin plots.
- Males do a little better in Maths
- Females do better in language subjects (reading & writing)



3.5 Analysing Distributions

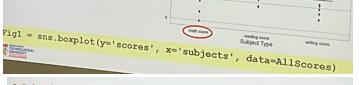
(3) Toughest subject

Which subject did the students do worst in?

Solution:

Use box plot.

Maths seems to be the most difficult subject for the average student.



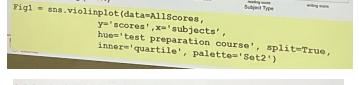
3.5 Analysing Distributions

(4) Test Preparation Course

Did the prep course help?

Solution:

- Use violin plot [for each level]
- Completing test preparation course seems to help improve average scores in all subjects.
- Most in writing (+11%)



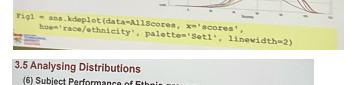
3.5 Analysing Distributions

(5) Performance of Ethnic groups

Did performances vary among the 5 groups?

Solution:

- Use KDE plot.
- Groups D and E did the best.
- Group A did the worst (82%)



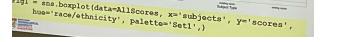
3.5 Analysing Distributions

(6) Subject Performance of Ethnic groups

Did subject performance vary among the 5 groups?

Solution:

- Use box plot.
- Group A did the worst in all three subjects.
- Group E did the best in all three subjects (especially Maths)



3.5 Analysing Distributions

(7) Parent education on student scores

Did parental education level influence student's score?

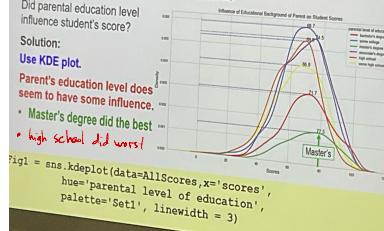
Solution:

Use KDE plot.

Parent's education level does seem to have some influence.

Master's degree did the best

high school did worst!



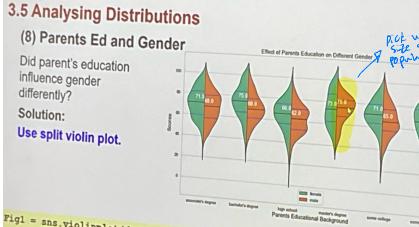
3.5 Analysing Distributions

(8) Parents Ed and Gender

Did parent's education influence gender differently?

Solution:

Use split violin plot.



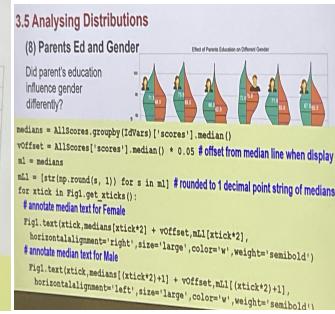
3.5 Analysing Distributions

(8) Parents Ed and Gender

Did parent's education influence gender differently?

Solution:

Use split violin plot.



4.1 Weber's Law**(To be Done Before Week #6 and #7)**

- (1) Using Powerpoint (*or any application where you can create grey level squares with any intensity values*), create eight reasonable-sized and equally-spaced squares with incrementing Red, Green, Blue (R, G, B) intensity values as shown in Figure 1.
- (2) With the eight grey level squares displayed within a white background, state which two squares have intensity levels that are most difficult to tell apart.
- (3) Now change your background to totally black, that is, the grey value of (0,0,0). Now see if the same two squares that you had problem distinguishing are still the same ones. If not, which two squares do you now find most difficult to distinguish? Can you explain why the background intensity can influence your ability to discriminate the grey value intensity?
- if change background to black then 6 and 7*
- (4) If you only had these eight discrete grey values to encode a set of ordinal values from your dataset, what background would you use on your chart to maximise your available distinguishable discrete luminance (intensity) attribute.
- discriminate value of #5 / (112, 112, 112)*
- (5) How would you change the given set of discrete grey values to improve the discriminability of these eight intensity attributes if you would like to use a white background?

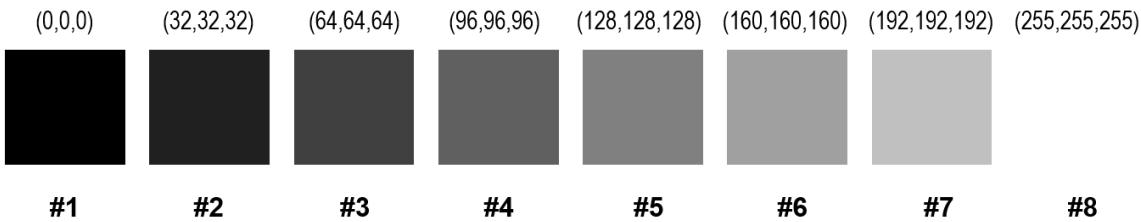


Figure 1 – Eight grey valued squares with equal intensity increments. The white square (#8) cannot be seen as it has the same intensity value as the white background.

4.2 Effective Visual Story

- (1) Table 1 shows the average measured levels of a fictitious hormone called *Vitalis* in a population based on age group, gender and BMI. The visual story you want to communicate is that the **only group** with increasing levels of Vitalis as they **age** are **females** with **BMI < 25**.

Design an appropriate chart to make this story **stand out visually** while providing all the information shown in the table.

Levels of the hormone Vitalis			
Body Mass Index (BMI)	Males		Females
	Under 60 years	60 years or over	Under 60 years
Under 25	255	230	380
25 or over	440	325	720

Table 1 – Average Vitalis levels measured in the sampled population under the various categories.

Note: The Excel file “*Hormone Vitalis – Excel.xlsx*” with Table 1 is given to you. You can generate your chart in Excel or use “*Hormone Vitalis – CSV.csv*” and do the plot in Matplotlib.

4.1 Weber's Law

(3) Which two squares are most difficult to distinguish?

(0,0)	(32,32,32)	(64,64,64)	(96,96,96)	(128,128,128)	(160,160,160)	(192,192,192)	(255,255,255)
#1	#2	#3	#4	#5	#6	#7	#8

Why the difference?

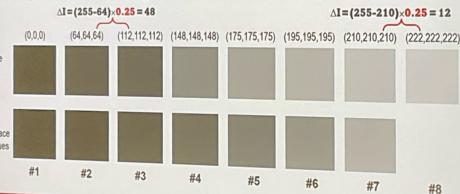
- The background Intensity changes what is considered a higher intensity stimulus as the Intensity of luminance stimulus we perceive is dependent on the amount of contrast between the background and foreground.

4.1 Weber's Law

$$\frac{\Delta I}{I} = k$$

(5) Change grey values to improve discriminability in white background?

- Compensated scale using a Weber's constant of **0.25** (first step is 64).



12

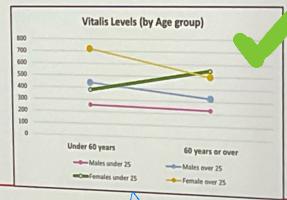
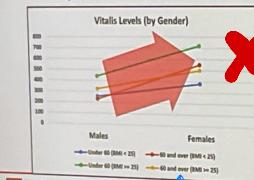
4.2 Effective Visual Story

(1) Design chart to make story stand out

(i.e. Females with BMI < 25) has increasing Vitalis as they get older

- This group should have a distinct visual pattern compared to the rest.

Body Mass Index (BMI)	Males		Females	
	Under 60 years	60 years or over	Under 60 years	60 years or over
Under 25	255	230	380	550
25 or over	440	325	720	500



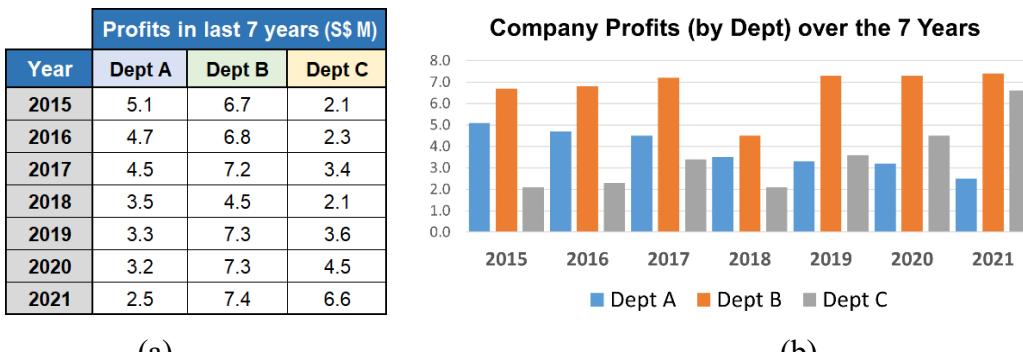
15

if male and female

stand out more

4.3 Applying Gestalt Principles

- (1) The table in Figure 2(a) shows the annual profit of each department (A, B and C) in an imaginary company. The data table, along with its bar chart shown in Figure 2(b) can be found in the Excel file “*Changing Profits – Excel.xlsx*”. An equivalent CSV version of the dataset is also provided (if you want to do the visualization in Python).



(a)

(b)

Figure 2 – (a) The data table showing how each department contributes to the company’s changing profit each year (2015 to 2021). (b) A bar chart showing each department’s annual profit.

- (2) **Presentation by Department Head** – As the head of department (Dept C), you are making a **Powerpoint** (PPT) presentation to your Chief Executive Officer (CEO) to request for more resources to maintain the current strong growth in your department. Apply appropriate Gestalt principles to improve the design of the presentation shown in Figure 2(b).

Note: You can use more than one chart in your PPT (if you wish), add animation, sequence overlay, etc to make build your case and communicate your purpose clearly and effectively.

In your presentation design, communicate the following:

- Firstly, you must show the annual profits for all three department over the 7 years but make sure the presentation is primarily focused on your department’s performance over these years.
- Show clearly the changing trend of the company’s overall profits from one year to the next.
- Highlight the really challenging year for the company and your department.
- Communicate clearly how well your department has been performing over the last 7 years, especially after that very challenging year for the company.
- Show how your Dept C’s percentage contribution to the company’s overall profit has been growing over the years.

- (3) **Presentation by CEO** – As the CEO of the company, you are making a **Powerpoint** (PPT) presentation to Board of Directors to re-distribute the limited resources in the company. You need to make a case that the poorest performing department will need to be restructured to re-allocated resources to the department with the strong growth trajectory. In your visualisation, apply appropriate Gestalt principles to visually communicate the following:

- Compare the performance of the three departments over the last 7 years with the goal of highlighting the profit trend of the poorest performing department (Dept A) and the department with the strong growth potential (Dept C).
- The changing percentage contributions of each department to the company’s overall profits.

4.3 Applying Gestalt Principles

(2) HOD (Dept C) to request for more resources to maintain the current strong growth

- **Proximity** - is used to group the departments together into the different years so it is easy to see the changing profits from year to year.
- **Similarity** - used in reverse to highlight Dept C's data as it is the only one in colour (Green connotes growth)
- **Connectedness** - An additional line chart with a secondary axis is used to show the company's overall profits.
- **Similarity** - In colour used to link right vertical axis to line plot.



Gestalt Principles used to improve the graph above?

4.3 Applying Gestalt Principles

(3a) CEO to Board – Redistribute resources, who is up and down...

- **Connectedness** - A line chart shows changing trend. A growing green colour is assigned to Dept C and a warning red is assigned to Dept A. Non-relevant Dept B de-emphasised with a grey colour.
- **Similarity & Proximity** - Annotation of relevant figures (using similar colour as the lines) are placed at near lines shows the declining and increasing profits of Dept B and C respectively over the last 7 years



Gestalt Principles used to improve the graph above?

4.3 Applying Gestalt Principles

(2c-d) Highlight challenging year and performance after 2018

- **Enclosure** - is used to highlight the challenging year of 2018. Global crisis is annotated over that year to explain why it is highlighted.
- **Connectedness** - A strong line is used to highlight Dept C's increasing profits.
- **Proximity** - Actual values are annotated at 2018 and at 2021 to show the amount of increase over the last four years.



Gestalt Principles used to improve the graph above?

4.3 Applying Gestalt Principles

(3b) CEO to Board – Percentage contributions

CEO



HOD(Dept C)



Contrast the two charts from the CEO and HOD (Dept C)

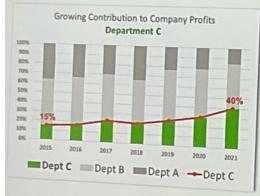
only HOD department

21

4.3 Applying Gestalt Principles

(2e) Highlight Dept C's percentage contribution over the years....

- **Similarity** - Percentage stacked bar chart to show the % contribution. Again, the same consistent green colour used to make Dept C's profit data stand out.
- **Connectedness** - Strong red line is highlight positive profit contribution trajectory of Dept C profits.
- **Proximity** - Actual % contributions are annotated on the bar charts (start and end) for further clarity of progress made so far.



Gestalt Principles used to improve the graph above?

4.3 Applying Gestalt Principles

(3b) CEO to Board – Percentage contributions

Changing Contribution to Company Profits



most say which Gestalt principles are involved!

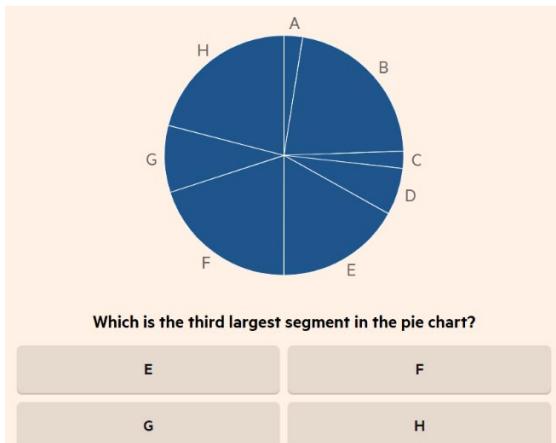
- A stacked bar chart is used to show the changing percentage contribution of each department.
- **Proximity** - A downward sloping red arrow near Dept A most recent profit data show recent declining profits.
- **Continuity & Similarity** - Colour similarity associates the data to each dept. Dept A & C are placed at the top & bottom ends as the principle of continuity with reference to horizontal axes will accentuate trends over the years.

What Gestalt Principles was used?

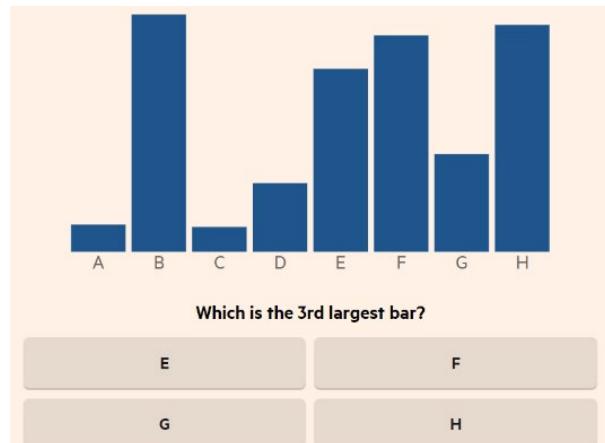
21

4.4 Making it Easy for Others

- (1) Financial Times provides an interesting online quiz entitled “*The science behind good charts*”, which was found at <https://ig.ft.com/science-of-charts/> (Note: subscription may be needed for access to FT website and link may be outdated). This online quiz helps you appreciate the feeling when you are at the receiving end of poorly designed data visualisations.



(a)



(b)

Figure 3 – (a) Pie Chart and (b) Bar chart taken from (<https://ig.ft.com/science-of-charts/>)

- (2) Figure 3 shows two questions taken from that FT quiz. Redesign the pie chart and the bar chart to make it “really-really” easy for someone taking the quiz to answer the question correctly.
- (3) Show what further improvements to the visualisation you would incorporate if your purpose is to get the audience to focus **specifically** and **precisely** on that third largest value in that dataset.
- (4) Using the Excel file “*FT Pie and Bar – Excel.xlsx*” or the CSV file “*FT Pie and Bar – CSV.csv*” provided, implement your improved version of the pie and bar charts in parts (2) and (3).
- (5) Describe what Gestalt principles were employed in improving your visualization.
- (6) Which of the two types of charts would you use if your purpose is to show how close entity **F** is to the next two highest rivals?
- (7) Are the exact differences between entity **F** and its top two rival entities **B** and **H** easy to visually estimate? If not, why is this so?
- (8) How can you make the differences between entity **F** and its top two rivals more apparent by:
- using the Gestalt principle of continuity.
 - using the Gestalt principle of enclosure.
- (9) Implement these improvements on the chart you have selected.

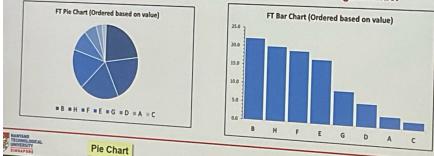
Note: You do not have to develop codes to implement every one of your improvement ideas. You can generate the basic pie or bar chart (using Excel or Python) and use appropriate applications with drawing tools (e.g. in Powerpoint, Word, Paint, etc) to add the additional improvement visuals into your base chart (exported out as an image file). You should spend more time thinking about the design instead.

4.4 Making It Easy for Others

(2) Redesign to make it really easy to get it right

What Gestalt principle was employed?

- The Gestalt principle of **simplicity** can be applied by sorting the pie segments and bars in a largest to smallest order to help the viewer see which is the 3rd largest value.

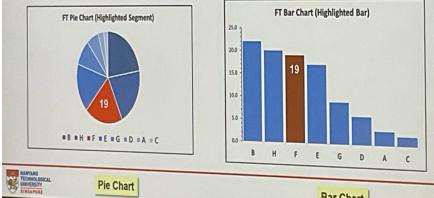


4.4 Making It Easy for Others

(3-4-5) How to get audience to focus on the 3rd largest value?

What Gestalt principle was employed?

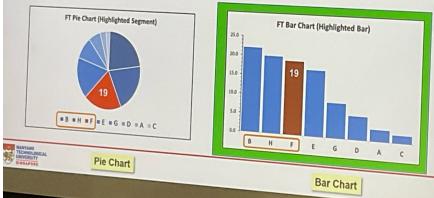
- The principle of **similarity (dissimilarity)** is used when highlighting (with a distinctively different colour) the particular third largest data value.



4.4 Making It Easy for Others

(6) Which chart to show how close entity F is to the next two highest rivals?

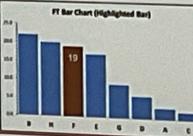
- A bar chart should be used if we want to compare data values more precisely.



4.4 Making It Easy for Others

(7) Are the exact differences between entity F and its top two rivals easy to visually estimate? Why?

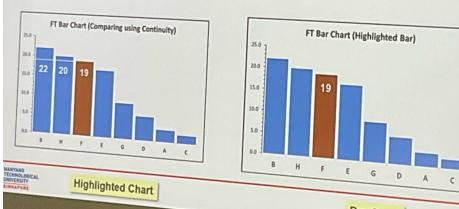
- Not easy to estimate the difference between F and B & H. The length disparity we are trying to visualise is much smaller than the overall length of the bar F. **Weber's Law** states that the increment must be larger if the intensity of the stimulus is large.



4.4 Making It Easy for Others

(8) Make differences betw. F and B & H apparent using **continuity**?

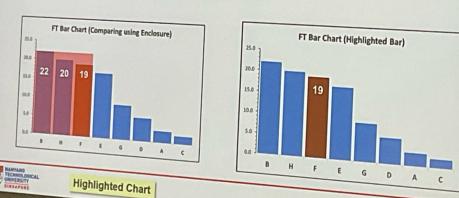
- Principle of continuity uses a white line to draw the reference height of entity F into the top two rivals.



4.4 Making It Easy for Others

(8) Make differences betw. F and B & H apparent using **enclosure**?

- The principle of enclosure uses the boundary of the enclosure to make the differences of the bar heights more visible.



4.5 Selecting Colour Palettes

- (1) Describe what are categorical (qualitative), sequential and diverging colour palettes; and under what condition would it be appropriate to use each of them. Using the Python Seaborn functions, create an example of each one of these palettes with 9 discrete and distinctive colours.
- (2) **Children Growth Rate** - The CSV file “*Children Growth Rate – CSV.csv*” charts the average yearly growth of boys and girls. There are interesting observations about the relative heights of boys and girls as they age from birth (0 years) to when they are 20 years old. Using **two stacked bars**, similar to that shown in Figure 4, and an **appropriate colour palette**, create a visualisation that will facilitate the comparison of growth rates between the gender. What did you observe?

Note: The children average growth rate dataset was modified from data obtained from *Disabled World* website at <https://www.disabled-world.com/calculators-charts/height-weight-teens.php>.

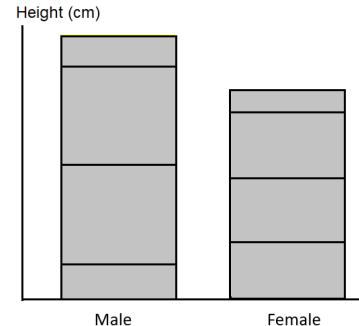


Figure 4 – A Sample Stacked Bar Chart

4.6 Less Colour but More Clarity

- (1) Figure 5 shows a data visualisation that presents the results from a survey regarding the work-life balance of various professional groups. List down what you think are poorly designed aspects of the visualisation that compromised the **accessibility** of the information presented.

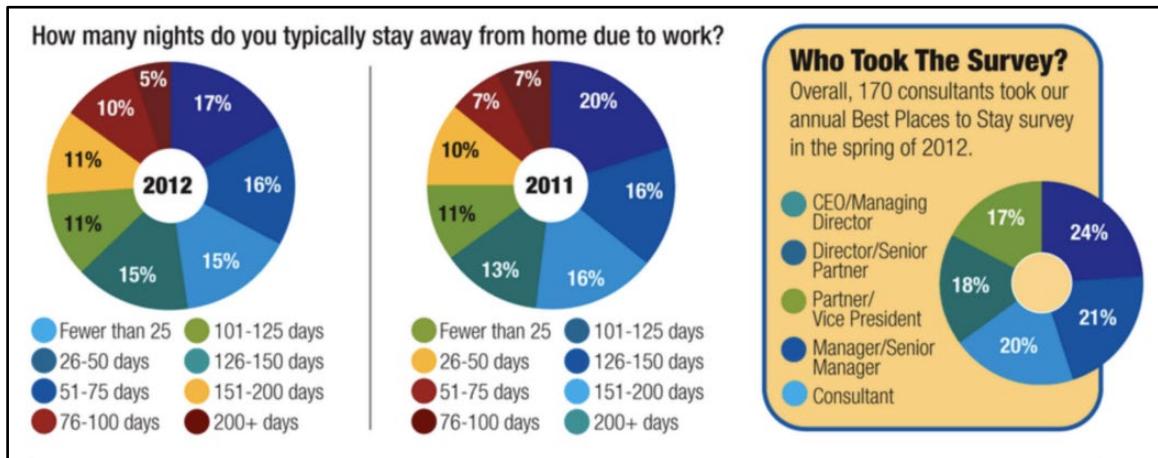
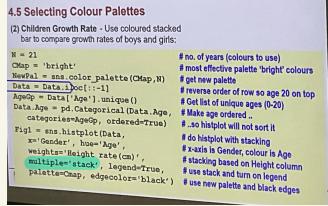
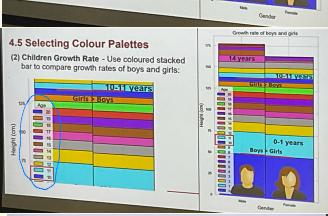
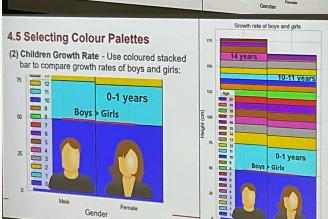
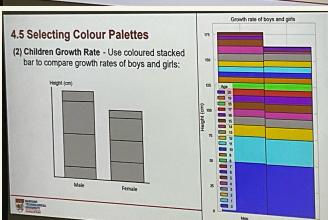
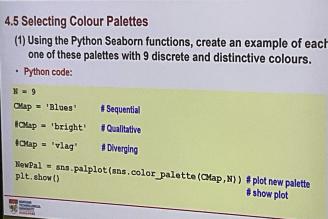
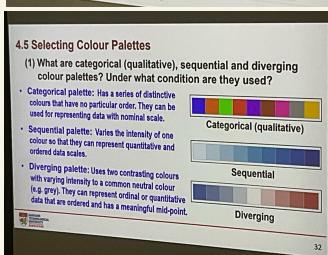
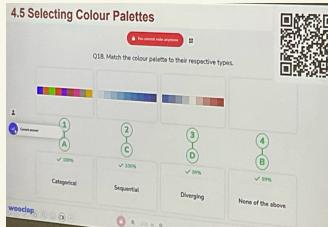


Figure 5 – A poorly designed visual that was featured in <https://towardsdatascience.com/color-in-data-visualization-less-how-more-why-348514a3c4d8>.

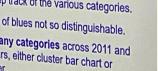
- (2) **No constrain redesign** – You are given no constraints in terms of chart type, number of charts, number of colours, use of annotations, etc. Redesign an improved version of the visual shown in Figure 5. Your new visualisation design must incorporate and show all the data and legend information present in Figure 5.
- (3) **One-colour only redesign** – Complete the redesign as in part (2) with no constraints except that now you can only use one colour (this one colour excludes the black for your text and borders and the white for your background).

Note: Those interested to create the visual can use the data in “*Away at Work – Excel.xlsx*” file.

**4.6 Less Colour but More Clarity**

(1) What are the poorly designed aspects of this visual?

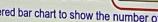
1. Colour assignment - for the two pie charts (2011 and 2012) should be consistent. The same category must have the same colour in both cases or else there is a frequent need to cross reference the colour legend.
2. Too many colours - making it difficult to keep track of the various categories.
3. Colours is not very distinct - three shades of blues not so distinguishable.
4. Pie charts - inappropriate to compare so many categories across 2011 and 2012. To see the change over these two years, either cluster bar chart or percentage stacked bar chart would be better.

*Solution***4.6 Less Colour but More Clarity**

(2) No Constraint Redesign

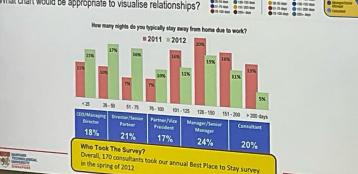
What chart would be appropriate to visualise relationships?

1. **The Two Pie Charts** - We could use a clustered bar chart to show the number of days distribution for 2011 and 2012.
2. **Clustered Bar chart** - makes comparison between the two years easy and is effective in telling the story that work-life balance is improving.
3. **Colour** - Can use the green colour for 2012 to convey this feeling of improvement.
4. **Survey Info** - No need for multiple colours to show survey participant distribution.
5. **Stacked % Bar chart** - The spread is even enough to allow us to put the text within the proportion box of the stacked percentage bar chart.

**4.6 Less Colour but More Clarity**

(2) No Constraint Redesign

What chart would be appropriate to visualise relationships?

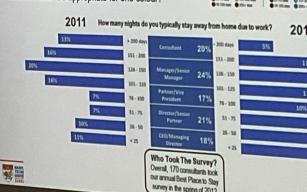


40

4.6 Less Colour but More Clarity

(3) One-colour Redesign

What chart would be appropriate for one colour?

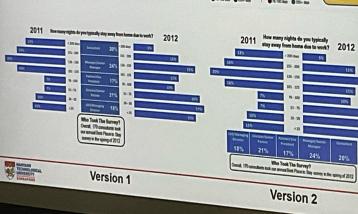


41

4.6 Less Colour but More Clarity

(3) One-colour Redesign

Version 1 or Version 2 do you prefer?



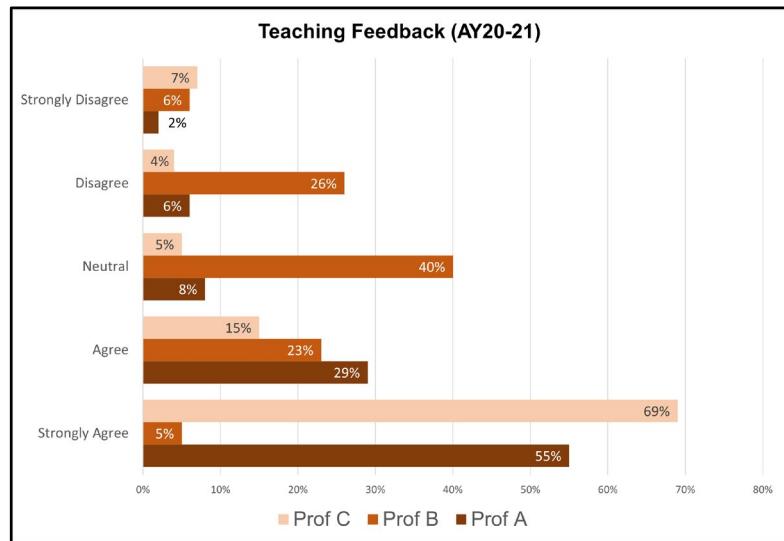
45

4.7 Critical Eye and Creative Fix

- (4) **Teaching Feedback Performance.** Figure 6(b) shows a fictitious clustered bar chart. It was created from the data table in Figure 6(a) of collated questionnaire responses tabulated as a percentage (%) of all responses received for each Professor **A**, **B** and **C** during a recent student teaching feedback survey. The survey was done based on a 5-point Likert scale, with the most favourable response stated as “Strongly Agree” and the least as “Strongly Disagree”.

Likert Rating	Prof A	Prof B	Prof C
Strongly Agree	55%	5%	69%
Agree	29%	23%	15%
Neutral	8%	40%	5%
Disagree	6%	26%	4%
Strongly Disagree	2%	6%	7%

(a)



(b)

Figure 6 – (a) Data table of teaching feedback score percentages over a 5-point Likert scale for three professors. (b) A clustered bar chart comparing the teaching performances of the 3 professors.

- (5) Is the data visualisation shown in Figure 6(b) effective in comparing and communicating the relative teaching performances of the three Profs **A**, **B** and **C** based on their respective distribution of 5-point Likert-based responses? Your comments should touch on the following:
- The choice of chart for the purposes stated above.
 - The choice of colours used in the chart.
 - Relevant Gestalt principles that have or have not been exploited to make the visualisation more effective.
- (6) Carefully consider the nature of the data and the intended purposes of the visualisation, then design an improved version. Implement your new design using the data from the Excel file “Teaching Feedback - Excel.xlsx” or the CSV file “Teaching Feedback – CSV.csv” provided.
- (7) Describe the rationale for your improved visualisation design, which should cover the following:
- Reasons for your choice of the particular chart type.
 - Choice of colours or colour palette used in the chart.
 - Relevant Gestalt principles you have employed in your improved design.
- (8) What can you infer when comparing the relative teaching feedback distributions of Profs **A**, **B** and **C**?

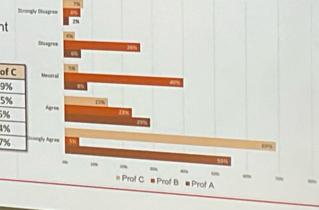
4.7 Critical Eye and Creative Fix

(1) Teaching Feedback Performance

- Data tabulated as a percentage of all responses received.
- Survey done based on a 5-point Likert scale.

Likert Rating	Prof A	Prof B	Prof C
Strongly Agree	55%	5%	69%
Agree	29%	23%	15%
Neutral	8%	40%	5%
Disagree	6%	26%	4%
Strongly Disagree	2%	6%	7%

Teaching Feedback (AY20-21)



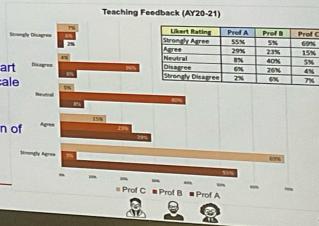
46

4.7 Critical Eye and Creative Fix

(2) Effectiveness in comparing teaching performance?

- Choice of chart?
- Choice of colour?
- Gestalt principles?

- a) **Choice of chart - cluster bar chart is inappropriate as the Likert scale distribution is a percentage.**
Stacked percentage bar chart better. It will show the distribution of the 5 categories for each Prof, making the comparison much easier.



Should be stacked
vs stacked

47

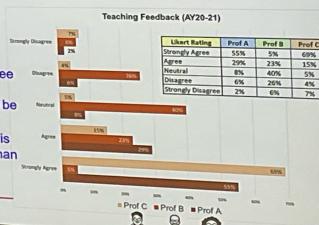
4.7 Critical Eye and Creative Fix

(2) Effectiveness in comparing teaching performance?

- Choice of chart?
- Choice of colour?
- Gestalt principles?

- b) **Choice of colour - Since the three Pros are nominal categories, a more distinctive palette would be better.**

The gradation of brown colours is closer to a sequential palette than a qualitative palette.



48

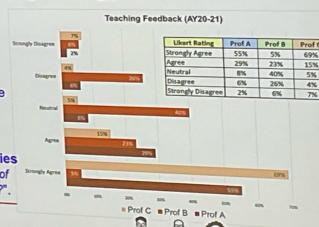
4.7 Critical Eye and Creative Fix

(2) Effectiveness in comparing teaching performance?

- Choice of chart?
- Choice of colour?
- Gestalt principles?

- c) **Gestalt principles - the use of proximity groups the Likert scale categories. Easier to compare scores in each category.**

But, proximity also makes it harder to view across categories to ask questions like "how much of the scores are positive, i.e. SA and A?"



49

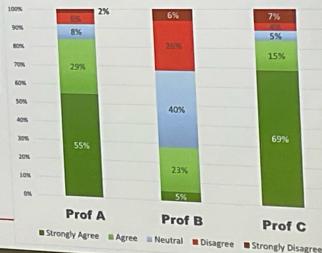
4.7 Critical Eye and Creative Fix

(3) Possible New Design

Likert Rating	Prof A	Prof B	Prof C
Strongly Agree	55%	5%	69%
Agree	29%	23%	15%
Neutral	8%	40%	5%
Disagree	6%	26%	4%
Strongly Disagree	2%	6%	7%



Teaching Feedback (AY20-21)



Teaching Feedback (AY20-21)

4.7 Critical Eye and Creative Fix

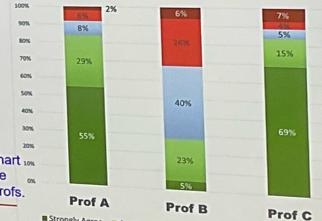
(4) Rationale for Design

- Choice of chart?
- Choice of colour?
- Gestalt principles?

- a) **Choice of chart - Percentage stacked bar chart is more appropriate as the Likert scale rating distribution is given as a percentage.**
Uniform height enforced by chart makes it easy to compare the rating distribution across the Pros.

Teaching Feedback (AY20-21)

Teaching Feedback (AY20-21)



Teaching Feedback (AY20-21)

4.7 Critical Eye and Creative Fix

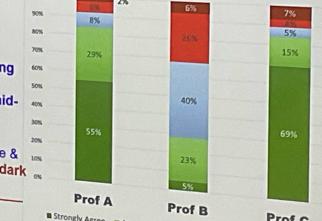
(4) Rationale for Design

- Choice of chart?
- Choice of colour?
- Gestalt principles?

- b) **Choice of colour - Use diverging palette as Likert scores have negative, positive and neutral mid-point.**
Colour match the negative and positive emotions. Green for +ve & reddish for the -ve, with a very dark shade of red for SD

Teaching Feedback (AY20-21)

Teaching Feedback (AY20-21)



Teaching Feedback (AY20-21)

4.7 Critical Eye and Creative Fix

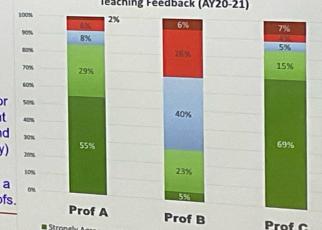
(4) Rationale for Design

- Choice of chart?
- Choice of colour?
- Gestalt principles?

- c) **Gestalt principles - principle of similarity is used when colour for each Likert category is consistent across the three stacked bars and they are made distinct (dissimilarity) from other categories.**
This way, it is easier to compare a given rating across the three Pros.

Teaching Feedback (AY20-21)

Teaching Feedback (AY20-21)



Teaching Feedback (AY20-21)

4.7 Critical Eye and Creative Fix

(5) What can you say about the Prof A, B and C?

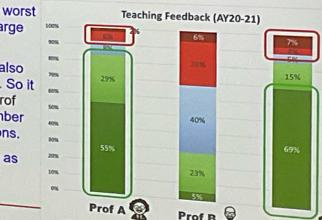
Prof B - is clearly the one with the worst teaching performance due to the large % of SD and D ratings.

Prof C - has largest % of SA. But also more % of SD compared to Prof A. So it seems many students really like Prof C's teaching but a reasonable number shared very strong negative opinions.

Prof A - has same 84% of SA & A as Prof C but fewer students strongly dislike his or her style of teaching.

Teaching Feedback (AY20-21)

Teaching Feedback (AY20-21)



Teaching Feedback (AY20-21)

4.8 Psychological Principles for Effective Graphics

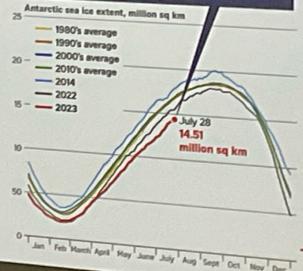
(1) Melting Antarctica Sea Ice

Uncharted territory

Every winter, Antarctica is locked in by a vast ring of sea ice. For decades, the area of winter sea ice had been measured every day, but in recent years, the area of that ice has started to shrink alarmingly and scientists fear climate change is to blame.



The sea ice extent for this time of year is by far the lowest, since satellite measurements started more than 40 years ago and has become yet another indicator of the record-breaking extremes happening around the globe.



Critique example

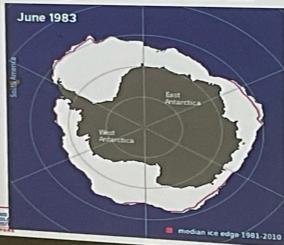


NOTE: US National Snow and Ice Data Center's (NSIDC) scientists use the 1981 to 2010 average of sea ice concentration and extent to have a consistent baseline for comparison to today's fluctuating conditions. The 30 years observed in this baseline provide enough data to even out short-term variability.

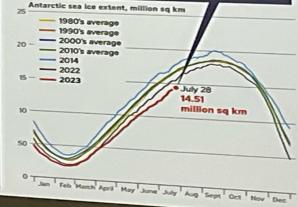
Sources: NSIDC, ARCTIC DATA ARCHIVE SYSTEM
STRATIS TIMES GRAPHICS

4.8 Psychological Principles for Effective Graphics

(1) Melting Antarctica Sea Ice



The sea ice extent for this time of year is by far the lowest since satellite measurements started more than 40 years ago and has become yet another indicator of the record-breaking extremes happening around the globe.



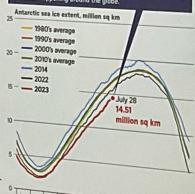
4.8 Psychological Principles for Effective Graphics

(1) Melting Antarctica Sea Ice

- a) What is the main message the line chart is attempting to communicate?

Main message - the sea ice extent for this time of the year is the lowest since the last 40 years of satellite measurements.

This is to emphasize the record-breaking extremes in climate change that is happening around the world.



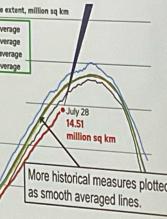
4.8 Psychological Principles for Effective Graphics

(1) Melting Antarctica Sea Ice

- b) Identify psychological principles and describe how it was used to make visual more effective.

Goldilocks Principle & Principle of Capacity Limitation

1980's to 2010's data were averaged and plotted as only four line plots to avoid clutter. This makes plots of recent changes in the melting sea ice extent more clearly observable.

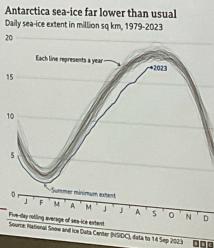


4.8 Psychological Principles for Effective Graphics

(1) Melting Antarctica Sea Ice

Data Visualisation by BBC News

What principle was employed?



4.8 Psychological Principles for Effective Graphics

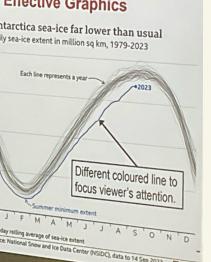
(1) Melting Antarctica Sea Ice

Data Visualisation by BBC News

What principle was employed?

Principle of Salience

The 2023 line plot (most recent year) was drawn in blue to focus user's attention to the latest sea ice extent measurements. All other lines in inconspicuous light grey



4.8 Psychological Principles for Effective Graphics

(1) **Melting Antarctica Sea Ice.** Study the line chart in Figure 7, which shows the data of the Antarctica sea ice extents from 1980 to 2023.

- What do you think is the main message the line chart is attempting to communicate to the reader?
- Identify which of Stephen Kosslyn's eight psychological principles of effective graphics were employed by The Straits Times to create the infographics shown. You may list more than one. (see Lecture notes chapter 7). For each of the principles that you have identified, describe how it was used to convey the intended message more effectively.

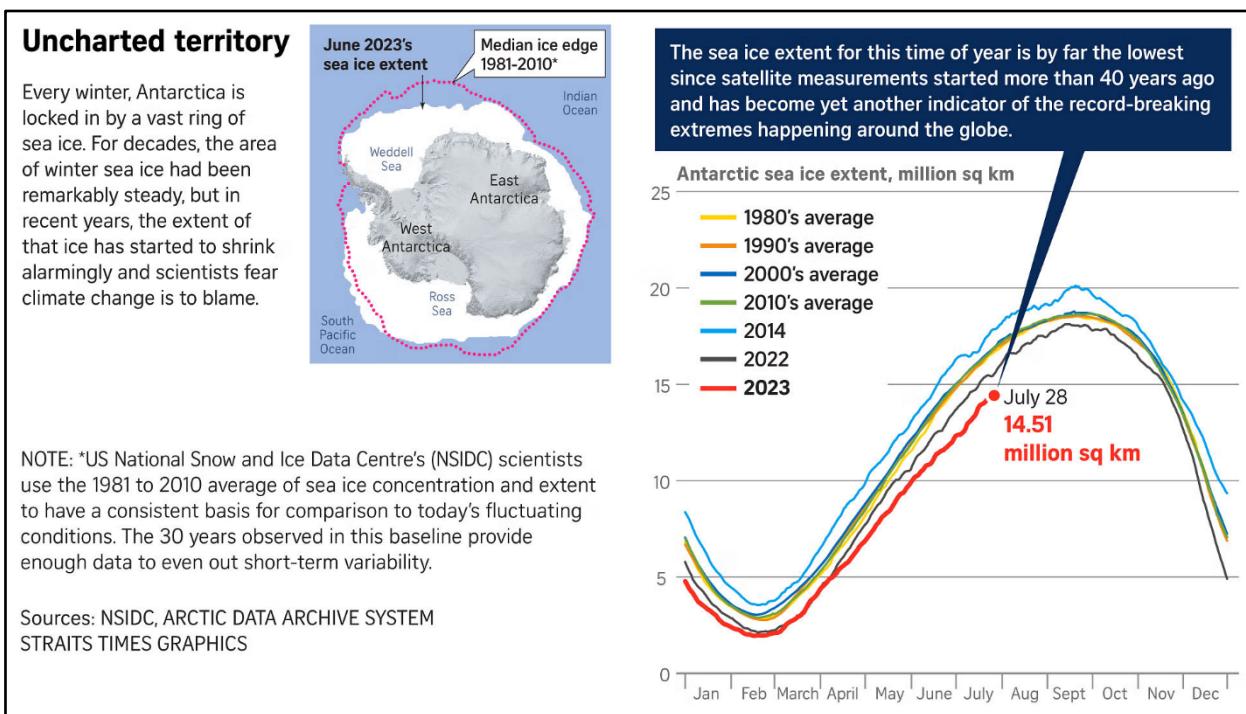


Figure 7 – The infographics illustrating the sea ice extent variations around the Antarctica from 1980 to 2023. This data visualisation was taken from an article in The Straits Times entitled “Antarctica’s record low sea ice worries climate scientist”, which was published on 30 July 2023. Link to article: <https://www.straitstimes.com/world/antarctica-s-record-low-sea-ice-worries-climate-scientists>.

part 2

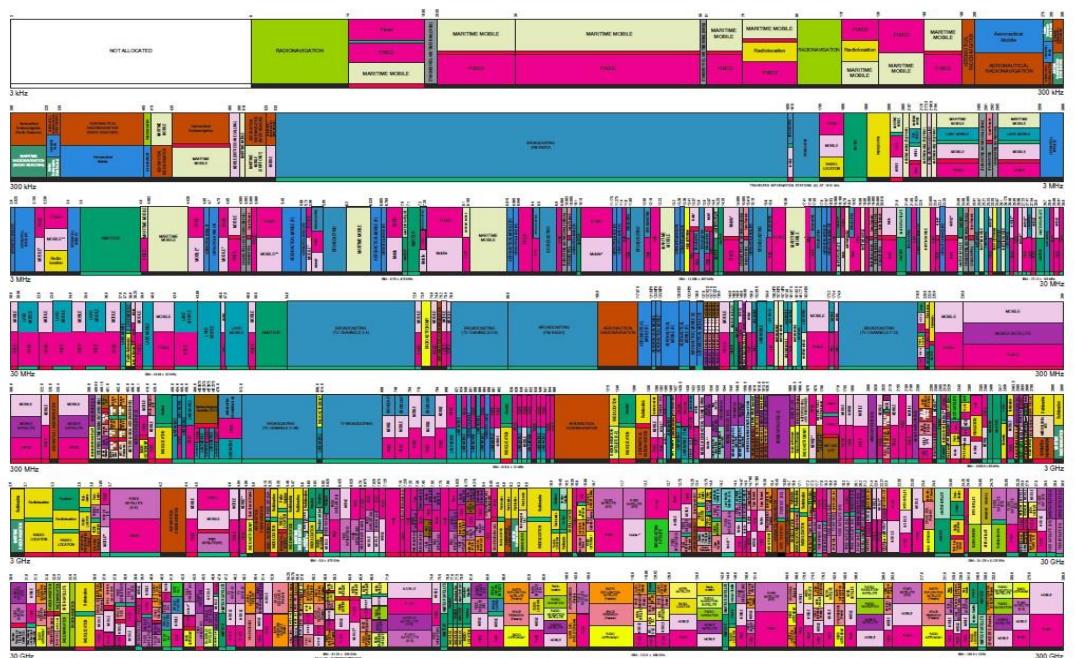
SC 4024 Tutorial 5: Visualization Principles and Interactions

1. Color Encoding for Radio Spectrum Allocation

The following figure shows the visualization for radio spectrum allocation in the United States, which is collected from here:

<https://www.ntia.gov/files/ntia/publications/2003-allochrt.pdf>

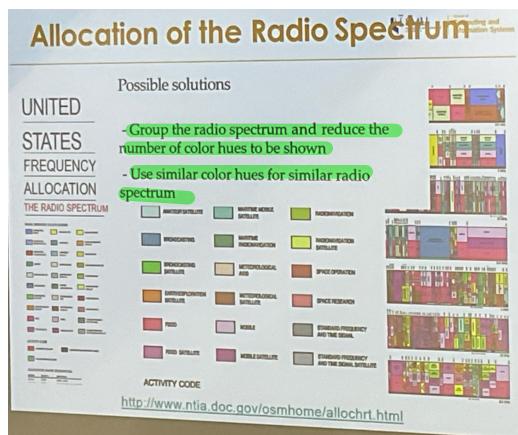
UNITED STATES FREQUENCY ALLOCATIONS THE RADIO SPECTRUM



1.1 By referring to the visualization principles and best practices we have learned in this course, please **point out the major issues of the color encoding design** of the above visualization.

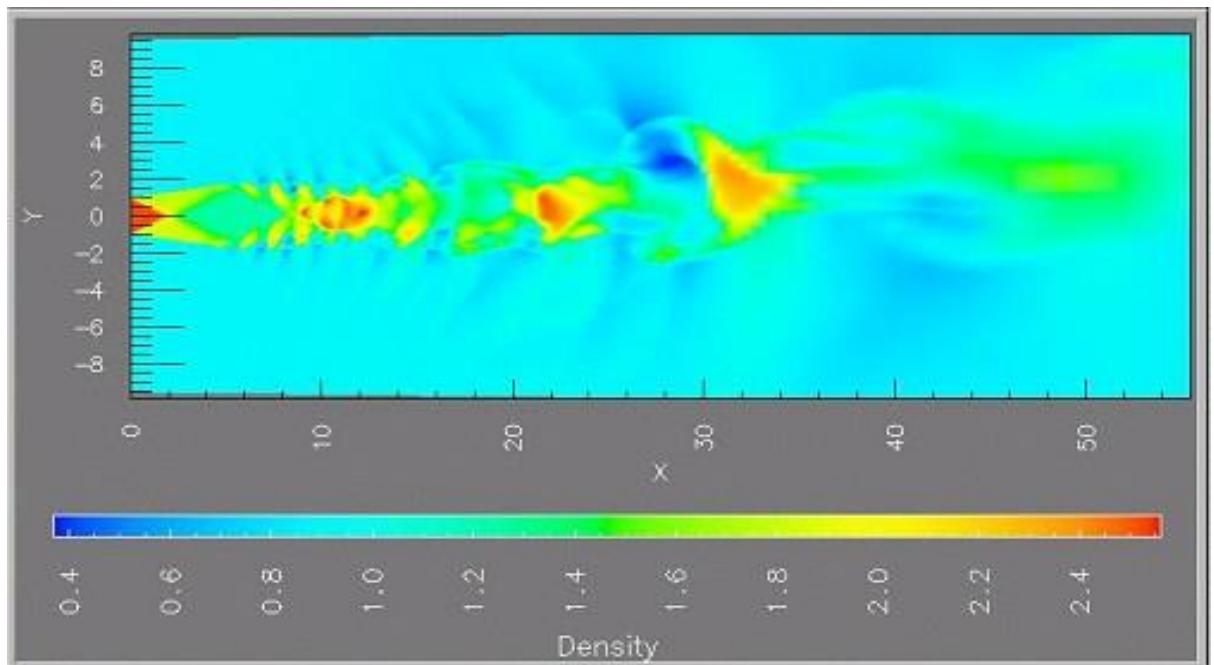
1.2 Suppose you are the visualization designer, do you have **any suggestions to address the above issues in terms of color encoding design?**

- 1.1)
- too many color, can't remember the mapping
 - Some color will be too similar (perception issue)



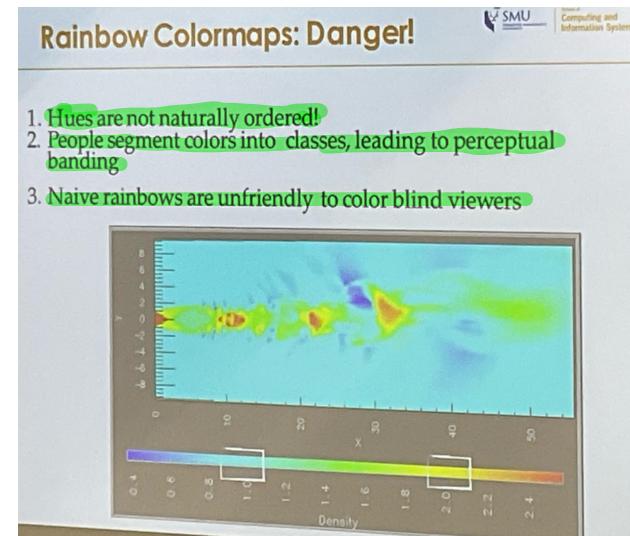
2. Rainbow Colormaps

A rainbow colormap is based on the order of colors in the spectrum of visible light—the same colors that appear in a rainbow. Below is an example of rainbow colormap, which is used to display the salt density of the sea in a specific area.



2.1 Does the above rainbow colormap look good to you? Why or why not? You are required to provide the detailed justifications.

2.2 If there are issues in the colormap of the above visualization, how should we improve it?



2.2 - Use seq or diverging colormap

SC 4024 Tutorial 6: Exploratory Data Analysis

AY2024-2025 Semester 1

1. Fundamental Concepts of Correlation Analysis

For the following Pearson correlation coefficients (r):

0.32, -0.76, 0.13, -0.25, 0.04, -0.03, 0.01

a) Which is the strongest correlation? - 0.76

b) Which is the weakest correlation? 0.01

2. Exploratory Data Analysis

During the lecture of Chapter 9.1, we take the course scores of students as an example to illustrate how we can conduct exploratory data analysis via data visualization and statistical analysis. Now you are provided with a dataset file called "students_performance.csv", and you are asked to conduct exploratory data analysis on it. Suppose you are interested in **checking if there is a difference between the math scores and reading scores of students from group A**, please answer the following questions:

- as long as
1 of attr not follow
normal distribution

Should not use
Z test*
- 2.1 By referring to what we have learned about exploratory data analysis and correlation analysis, briefly describe **what kind of steps** you will take to answer such a question, and what kind of **visualizations** and **statistical tests** you will use.
*- Check dataset size
- Parametric / non-parametric test
- Use Z-test*
 - 2.2 What are your **null hypothesis (H_0)** and **alternate hypothesis (H_1)** for such an analysis?
*Null → there is no diff in math & write
alte → there is diff in math & write*
 - 2.3 By following the steps in 2.1, you are asked to write Python code to create appropriate visualizations and conduct the corresponding relevant statistic tests. Note: You are requested to use the Python packages like **scipy.stats** and **statsmodels.stats**.

P < 0.05 → reject null hypothesis

P ≥ 0.05 fail to reject null hypothesis

2.4 Suppose you are asked to check there is a difference between the math scores and reading scores of **students from group B or group C**, what kind of changes do you need to do in your procedures of exploratory data analysis in 2.3?

null : there is no diff between my score n reading score $D.17 > 0.05$
alter : there is a diff in math n writing score

SC 4024 Tutorial 7: Geospatial Data Visualization

AY2024-2025 Semester 1

1. Selecting Appropriate Geospatial Data Visualization

We have learned a series of geospatial data visualization approaches for different types of geospatial data. Suppose you are requested to visualize the following data, what kind of geospatial data visualization approaches will you use? Justify your design choice.

- a) The annual sales of each McDonald's at different locations of Singapore in 2023. *Point based → Symbol map* *Choropleth*
- b) The average temperature of different countries in the EU. *Choropleth*
- c) The aging population (age>65) of different states in the U.S. *Cartogram*

2. Symbol Map Implementation

You are provided with a CSV file called “sg_universities.csv” containing the detailed location and student size information of five universities in Singapore. You are requested to use **d3.js** (more specifically, **d3.geo**) package to create a symbol map to visualize the student size of all the provided universities on a map. Below are detailed requirements:

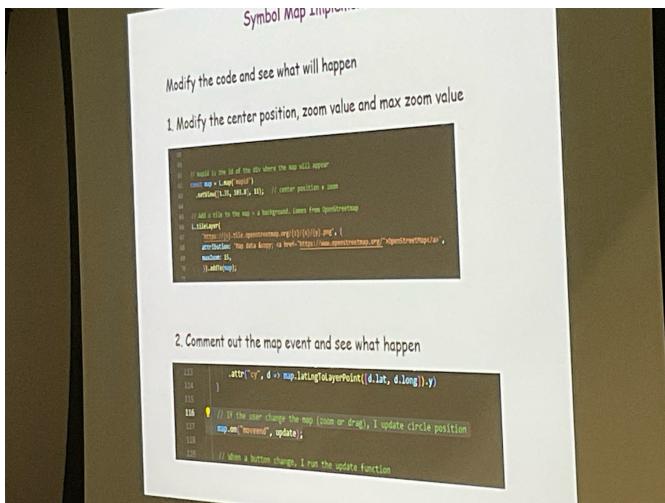
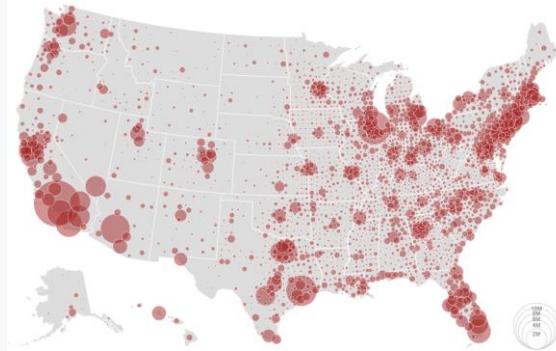
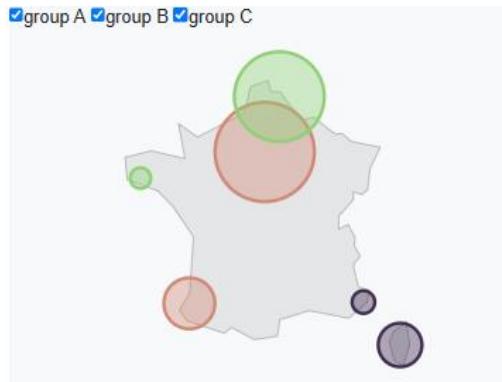
- a) **Visual encoding:** use a circle to represent each university with its radius encoding the student size (i.e., use bubble map to visualize the geospatial data), and use different fill color for different universities;
- b) **Interaction:** provide five checkbox to allow users to interactively select whether a university should be shown on the map or not;
- c) Animation: when a university is selected or un-selected, use animation to show or remove the glyph (i.e., circle);
- d) You are requested to use **D3 version 7** to implement the symbol map.

You are required to finish it before the tutorial on **Nov. 5th**.

References:

- [1] Bubblemap with filter button: https://d3-graph-gallery.com/graph/bubblemap_buttonControl.html

- [2] Bubblemap: <https://observablehq.com/@d3/bubble-map/2>





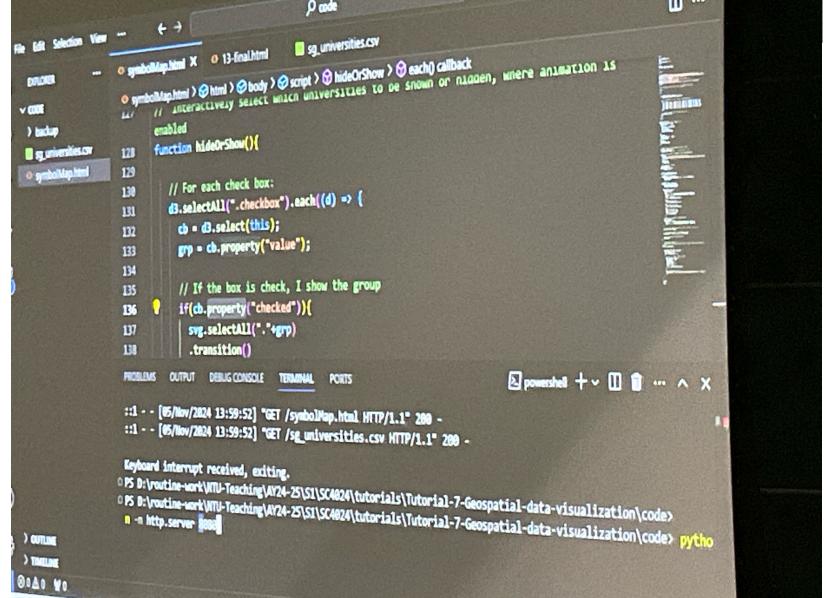
```
// Create data for circles
var markers = [
    {long: 1.39, lat: 43.34, group: "A", size: 10}, // Corse
    {long: 7.26, lat: 43.73, group: "A", size: 10}, // Ile de
    {long: 2.39, lat: 40.64, group: "B", size: 10}, // Poitou
    {long: 1.07, lat: 43.55, group: "C", size: 10}, // Bourgogne
    {long: 1.05, lat: 50.60, group: "C", size: 10}, // Lille
    {long: -3.0, lat: 40, group: "D", size: 10} // Melide
]

// Load external data and host
$.getJSON("https://raw.githubusercontent.com/bailey03/graph-gallery/master/JSONs/regions.json", function(data) {
    data.features = data.features.filter(function(feature) { return feature.properties.name !== null; });

    // Create a color scale
    var color = d3.scale.ordinal()
        .domain(["A", "B", "C", "D"])
        .range(["#4CAF50", "#FF9800", "#E91E63", "#9E9E9E"]);

    // Add a scale for bubble size
    var size = d3.scale.linear()
        .domain([1, 100]) // What's in the data
        .range([4, 10]); // Size in pixel

    // Draw the map
    var appender = g.selectAll("g")
        .selectAll("path")
        .data(data.features);
})
```



Tutorial 8

Graph Visualization

WANG Yong
College of Computing and Data Science
NTU

Force-directed Layout

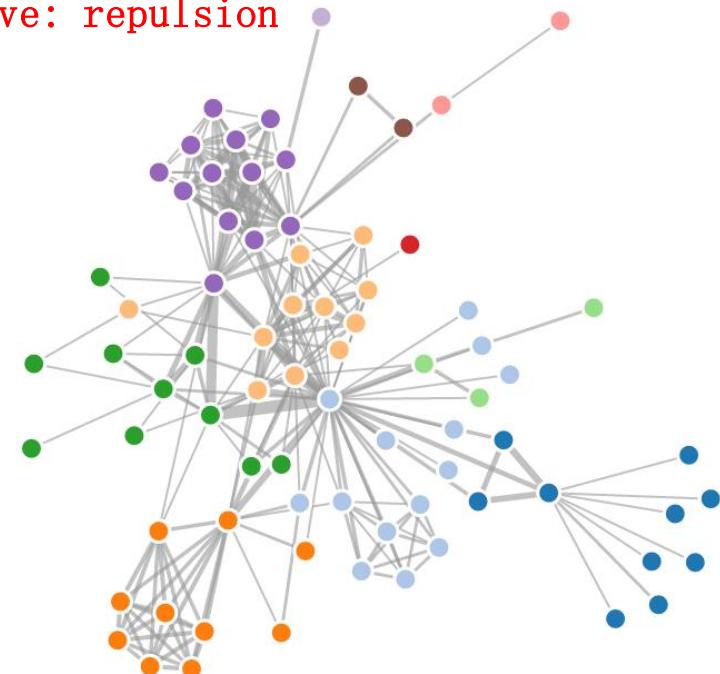
- There can be multiple forces between each pair of nodes.
- The layout is determined by all the forces.

```
var chargeForce = d3.forceManyBody()  
  .strength(0.1) //positive: attraction, negative: repulsion  
  .distanceMax(500)  
  .distanceMin(60);
```

```
// Creates a new circle collision force with  
// the specified radius
```

```
var collideForce = d3.forceCollide()  
  .strength(2)  
  .radius(20);
```

```
var simulation = d3.forceSimulation()  
  .force("link", d3.forceLink())  
  .force("charge", chargeForce)  
  .force('collide', collideForce)  
  .force("center", d3.forceCenter(width / 2, height / 2));
```



Note: we are using d3@v7; different versions of d3 can have slightly-different implementations!

Useful References

Official API:

<https://github.com/d3/d3/blob/main/API.md#forces-d3-force>

<https://github.com/d3/d3/blob/main/API.md#forces-d3-force>

<https://github.com/d3/d3-force/blob/v3.0.0/README.md#forceManyBody>

Examples:

<https://www.d3indepth.com/force-layout/>

<https://observablehq.com/@d3/force-directed-graph>

Your Task

1. Download the skeleton code
2. General task: finish the code to draw a force-directed graph

```
57     .attr("stroke-width", function(d) {  
58         // TO-DO: change the stroke-width as the square root of the "value" of the current link  
59         return 4;  
60     });
```

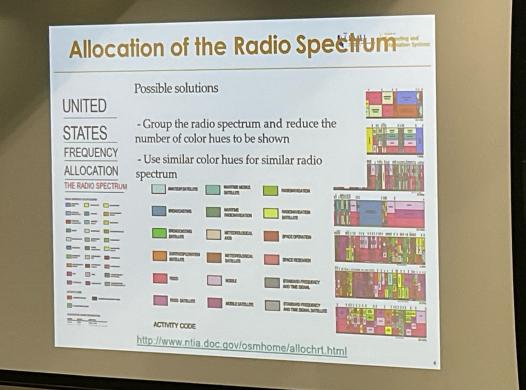
```
70     .attr("fill", function(d) {  
71         // TO-DO: revise the code to use different color to represent different groups  
72         return 'grey';  
73     })
```

Your Task

1. Download the skeleton code
2. General task: finish the code to draw a force-directed graph
3. Play with different parameter settings. Check the function of each parameter.

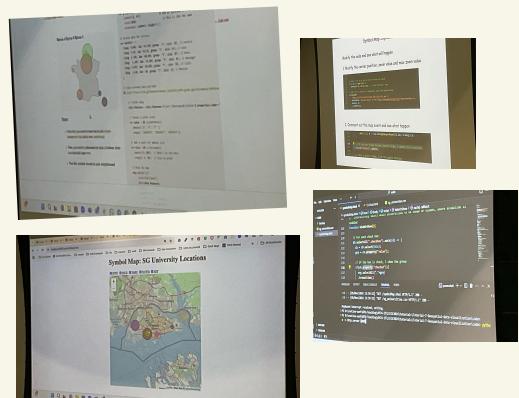
```
31 // TO-DO:  
32 //   1) Try to modify the values of strength, distanceMax and distanceMin, radius and see what will happen  
33 //   2) Try to use only "var chargeForce = d3.forceManyBody()" and see what will happen  
34 var chargeForce = d3.forceManyBody().strength(0.1).distanceMax(500).distanceMin(60);  
35 var collideForce = d3.forceCollide().strength(2).radius(20);  
36
```

THANK YOU~



tut 7
 1 a Symbol map
 b Choropleth map
 c Cartogram
 tut 6 Q3 part 2

TUT 7 Q2



1. Hues are not naturally ordered!
2. People segment colors into classes, leading to perceptual banding
3. Naive rainbows are unfriendly to color blind viewers



Strong Correlation is 0.76
 Weakest Correlation is 0.01
 → See absolute

Specifying scales

D3 scales are functions that map from a domain, to a range. (For instance, from the data domain, to the color of the chart.)
 Anomalous functions can have two parameters (our bound datum) and (the index of our datum).

```

    d = [0, 100, 200]
    s = d3.scale.linear()
        .domain(d)
        .range([0, 444])
    
```

Using a scale:
 scales
 .color("steelblue").domain(d).return(s))
 Manual specification:
 .color("steelblue").domain(d).range([0, 444])

Specifying scales

To position the dots, we must specify the x and y position attributes, but the process can be tedious and error prone for complex attributes.

```

    scatter
    .attr("cx", function(d){ return (380+1/sampleData.length)*d.x; })
    .attr("cy", function(d){ return 405-(d.y-2.574*(d.x/27.5)); })
    
```

selectAll().data().enter().append()

- Select all of our circles (currently we don't have any).
- Bind our data (in this case, 5 rows worth).
- Enter each new datum into our selection.
- Append a new DOM element. There are now 5 new elements, each with their own unique data.
- Set attributes with operators, using anonymous functions.

```

    var names = ["John", "Jane", "Mike", "Sarah", "David"]
    // Create a new circle for each name
    var circles = d3.selectAll("circle")
        .data(names)
        .enter()
        .append("circle")
        .attr("cx", function(d){ return (380+1/sampleData.length)*d; })
        .attr("cy", function(d){ return 405-(d.y-2.574*(d.x/27.5)); })
        .attr("r", 10)
        .attr("fill", "steelblue")
        .attr("stroke", "black")
        .attr("stroke-width", 1)
    
```

Questions & Experiments

- What is the major difference between Demo 4 and Demo 3 in terms of visualizing the same dataset as the same scatterplot?
- Can we use `d3.selectAll("circle").datum()` to pass the data associated with the circles in Demo 3?
- Try to change the circle attributes (e.g., fill-opacity, radius, stroke-width, cx, cy) and see what happens
- By revising the code of Demo 4, try to draw rectangles instead of circles!