

CSE 142 Assignment 1, Fall 2023

4 Questions, 100 pts, due: 23:59 pm, Oct 12th, 2023

Your name: Airi Kokuryo

Student ID: 2086695

Instruction

- Submit your assignments onto **Gradescope** by the due date. Upload a `zip` file containing:
 - (1) The saved/latest `.ipynb` file.
 - (2) All other materials to make your `.ipynb` file runnable.
- This is an **individual** assignment. All help from others (from the web, books other than text, or people other than the TA or instructor) must be clearly acknowledged.
- Most coding parts can be finished with only 1-2 lines of codes.
- Make sure you have installed required packages: `pandas`, `seaborn`, `matplotlib`

Objective

- **Task 1:** Review of **Probability and Linear Algebra**
- **Task 2:** Getting familiar with **Pandas** and **Seaborn/Matplotlib**

Question 1.1 (Conditional probability, 10 pts)

Assume that the conditional probability of an email (chosen uniformly and randomly from a set of emails) containing the word "payment", given that the email is a spam email, is 72%. Suppose that the conditional probability of an email being spam, given that it contains the word "payment", is 8%. Find the ratio of the probability that an email is spam to the probability that an email contains the word "payment".

Solution:

If you are not familiar with Latex, you may attach a figure/screen-shoot and display the code below.

```
In [ ]: from IPython.display import Image
        # Replace the figure name
        Image(filename='question11.jpg')
```

Out[]:

<Question 1.1>

probability of an email

- including the word "payment" ... $p(a)$
- is a spam email ... $p(b)$

from the passage,

$$p(a|b) = 0.72 = \frac{p(a \cap b)}{p(b)}$$

$$p(b|a) = 0.08 = \frac{p(a \cap b)}{p(a)}$$

$$\frac{p(b|a)}{p(a|b)} = \frac{p(b)}{p(a)} = \frac{0.72}{0.08} = \frac{9}{1}$$

$$\therefore p(b) : p(a) = \underline{9 : 1} //$$

Question 1.2 (Conditional probability, 10 pts)

There are two boxes. Box 1 contains three red and five white balls and box 2 contains two red and five white balls. A box is chosen at random $p(\text{box} = 1) = p(\text{box} = 2) = 0.5$ and a ball chosen at random from this box turns out to be red. What is the posterior probability that the red ball came from box 1?

Solution:

If you are not familiar with Latex, you may attach a figure/screen-shoot and display the code below.

```
In [ ]: from IPython.display import Image
# Replace the figure name
Image(filename='question12.jpg')
```

Out[]:

<Question 1.2>

$p(\text{box}=1) = p(\text{box}=2) = 0.5$ box 1: 3 red, 5 white
box 2: 2 red, 5 white

$p(\text{box}1|\text{red}) = ?$

$p(\text{red}|\text{box}1) = \frac{\frac{3}{8}}{\frac{1}{2}}$

$= \frac{3}{4}$

$\frac{3}{8} + \frac{2}{8} = \frac{21+16}{56} = \frac{37}{56}$

$p(\text{box}1|\text{red}) = \frac{p(\text{red}|\text{box}1) \cdot p(\text{box}1)}{p(\text{red})}$

$= \frac{\frac{3}{4} \cdot \frac{1}{2}}{\frac{37}{56}} = \frac{3}{8} \cdot \frac{56}{37} = \frac{21}{37}$

Question 1.3 (Gaussian & Poisson Distribution, 10 pts)

Part a) Gaussian Distribution

Let $X \sim N(0, 1)$ be a Gaussian random variable, which has the following probability density function:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Find $E(X)$ and show all the derivation steps.

Solution:

If you are not familiar with Latex, you may attach a figure/screen-shoot and display the code below.

```
In [ ]: from IPython.display import Image
# Replace the figure name
Image(filename='question13.jpg')
```

Out[]:

< Question 1.3 >

$$X \sim N(0, 1)$$

$$\underline{f(x)} = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\text{here, } \mu=0, \sigma=1$$

$$\underline{E(x)} = \int_{-\infty}^{\infty} x f(x) dx$$

$$= \int_{-\infty}^{\infty} \frac{x}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx$$

$$= \lim_{n \rightarrow \infty} \int_{-n}^n \frac{x}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx$$

$$= \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}} \int_{-n}^n \left\{ -\frac{1}{2} \exp\left(-\frac{x^2}{2}\right) \right\}' dx$$

$$= \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}} \cdot \left(-\frac{1}{2}\right) \left[\exp\left(-\frac{x^2}{2}\right) \right]_{-n}^n$$

$$= \lim_{n \rightarrow \infty} \left(-\frac{1}{2\sqrt{2\pi}} \right) \left[\exp\left(-\frac{n^2}{2}\right) - \exp\left(-\frac{(-n)^2}{2}\right) \right]$$

$$= \lim_{n \rightarrow \infty} \left(-\frac{1}{2\sqrt{2\pi}} \right) \cdot 0 = \underline{0}$$

Part b) **Poisson Distribution**

Let Y be a Poisson random variable, which has the following probability density function:

$$f(y; \lambda) = \Pr(Y=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

Where: 1) k is the number of occurrences ($k=0,1,2,\dots$)

2) e is Euler's number ($e=2.71828\dots$)

3) ! is the factorial function

Find $E(Y)$ and show all the derivation steps.

Solution:

If you are not familiar with Latex, you may attach a figure/screen-shoot and display the code below.

```
In [ ]: from IPython.display import Image
# Replace the figure name
Image(filename='question13b.jpg')
```

< Question 1.3 b >

$$f(y; \lambda) = \Pr(Y = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$E(Y) = \sum_{k=0}^{\infty} k \Pr(Y = k)$$

$$= \sum_{k=1}^{\infty} k \cdot \frac{\lambda^k e^{-\lambda}}{k!} dk$$

$$= \sum_{k=1}^{\infty} \frac{\lambda \cdot \lambda^{k-1} e^{-\lambda}}{(k-1)!} dk$$

$$= \lambda e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} dk$$

if put $k-1 = t$, $dk = dt$

$$= \lambda e^{-\lambda} \sum_{t=0}^{\infty} \frac{\lambda^t}{t!} dt$$

$$= \lambda e^{-\lambda} e^{\lambda} = \underline{\lambda}$$

Question 1.4 (Expectation and Variance, 10 pts)

Suppose that X_1, \dots, X_n are independent random variables with the same distribution.

(a --5 pts): Denote the mean of X_i as $\mathbf{E}[X_1]$, find the mean of

$$\frac{X_1 + \dots + X_n}{n}.$$

(b --5 pts): Denote the variance of X_i as $\text{Var}[X_1]$, find the variance of

$$\frac{X_1 + \dots + X_n}{n}.$$

Solution (a):

Solution (b):

If you are not familiar with Latex, you may attach a figure/screen-shoot and display the code below.

```
In [ ]: # Replace the figure name
from IPython.display import Image
Image(filename='question14.jpg')
```

Out[]:

< Question 1.4 >

(a)
$$\begin{aligned} \mathbf{E}[(X_1 + \dots + X_n)/n] &= \frac{1}{n} * (\mathbf{E}[X_1] + \dots + \mathbf{E}[X_n]) \\ &= \frac{1}{n} * n * \mathbf{E}[X_1] \\ &= \underline{\mathbf{E}[X_1]} \end{aligned}$$

(b)
$$\begin{aligned} \text{Var}[(X_1 + \dots + X_n)/n] &= (\frac{1}{n})^2 * (\text{Var}[X_1] + \dots + \text{Var}[X_n]) \\ &= (\frac{1}{n})^2 * (\text{Var}[X_1] * n) \\ &= \underline{\frac{\text{Var}[X_1]}{n}} \end{aligned}$$

Question 2 (Linear Algebra Review, 10 pts)

Find the (1) trace, (2) determinant, (3) matrix inverse, (4) eigenvalues & eigenvectors and (5) the eigenvalue decomposition of the following matrix (equal points for each).

$$A = \begin{bmatrix} 1 & 3 & 6 \\ 2 & 1 & 4 \\ 1 & 0 & 3 \end{bmatrix}$$

Solution:

If you are not familiar with Latex, you may attach a figure/screen-shoot and display the code below.


```
In [ ]: # Replace the figure name
from IPython.display import Image
Image(filename='question2.jpg')
```

Out []: <Question 2>

(1) trace

$$1 + 1 + 3 = 5$$

(2) determinant

$$\det \begin{bmatrix} 1 & 3 & 6 \\ 2 & 1 & 4 \\ 1 & 0 & 3 \end{bmatrix} = 1 \cdot \begin{vmatrix} 1 & 4 \\ 0 & 3 \end{vmatrix} + 3 \cdot \begin{vmatrix} 2 & 4 \\ 1 & 3 \end{vmatrix} + 6 \cdot \begin{vmatrix} 2 & 1 \\ 1 & 0 \end{vmatrix} \\ = (3-0) + 3(6-4) + 6(0-1) \\ = 3 + 6 - 6 = 3$$

(3) matrix inverse

$$A^{-1} = \frac{1}{|A|} A^T = \frac{1}{3} \begin{bmatrix} 1 & 2 & 1 \\ 3 & 1 & 0 \\ 6 & 4 & 3 \end{bmatrix}$$

$$\begin{vmatrix} 1 & 0 \\ 4 & 3 \end{vmatrix} = 3 \quad \begin{vmatrix} 3 & 0 \\ 6 & 3 \end{vmatrix} = 9 \quad \begin{vmatrix} 3 & 1 \\ 6 & 4 \end{vmatrix} = 6$$

$$\begin{vmatrix} 2 & 1 \\ 4 & 3 \end{vmatrix} = 2 \quad \begin{vmatrix} 1 & 1 \\ 6 & 3 \end{vmatrix} = -3 \quad \begin{vmatrix} 1 & 2 \\ 6 & 4 \end{vmatrix} = -8$$

$$\begin{vmatrix} 2 & 1 \\ 1 & 0 \end{vmatrix} = -1 \quad \begin{vmatrix} 1 & 1 \\ 3 & 0 \end{vmatrix} = -3 \quad \begin{vmatrix} 1 & 2 \\ 3 & 1 \end{vmatrix} = -5$$

$$\begin{bmatrix} 3 & 9 & 6 \\ 2 & -3 & -8 \\ -1 & -3 & -5 \end{bmatrix} * \begin{bmatrix} + & - & + \\ - & + & - \\ + & - & + \end{bmatrix} = \begin{bmatrix} 3 & -9 & 6 \\ -2 & -3 & 8 \\ -1 & 3 & -5 \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} 3 & -9 & 6 \\ -2 & -3 & 8 \\ -1 & 3 & -5 \end{bmatrix} * \frac{1}{3} = \begin{bmatrix} 1 & -3 & 2 \\ -\frac{2}{3} & -1 & \frac{8}{3} \\ -\frac{1}{3} & 1 & -\frac{5}{3} \end{bmatrix}$$

(4) eigen values & eigen vectors.

$$\lambda E - A = \lambda E - \begin{bmatrix} 1 & 3 & 6 \\ 2 & 1 & 4 \\ 1 & 0 & 3 \end{bmatrix} = \begin{bmatrix} \lambda-1 & -3 & -6 \\ -2 & \lambda-1 & -4 \\ -1 & 0 & \lambda-3 \end{bmatrix}$$

$$(\lambda-1)(\lambda-1)(\lambda-3) + (-3)(-4)(-1) + 0$$

$$-0 - (-3)(-2)(\lambda-3) - (-1)(\lambda-1)(-6)$$

$$= \lambda^3 - 2\lambda^2 + \lambda - 3\lambda^2 + 6\lambda - 3 - 12 - 6\lambda + 18 - 6\lambda + 6$$

$$= \lambda^3 - 5\lambda^2 + 7\lambda - 3 - 12\lambda + 12$$

$$= \lambda^3 - 5\lambda^2 - 5\lambda + 9 = (\lambda-1)(\lambda^2 - 4\lambda - 9)$$

$$\therefore \text{eigen values: } \lambda = 1, 2 \pm \sqrt{13}$$

(i) eigen vector of eigen value $\lambda = 1$

$$\lambda E - A = \begin{bmatrix} 0 & -3 & -6 \\ -2 & 0 & -4 \\ -1 & 0 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ -1 \end{bmatrix}$$

(ii) eigen vector of eigen value $\lambda = 2 + \sqrt{13}$

$$\lambda E - A = \begin{bmatrix} 1+\sqrt{13} & -3 & -6 \\ -2 & 1+\sqrt{13} & -4 \\ -1 & 0 & 1+\sqrt{13} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1+\sqrt{13} \\ 2 \\ 1 \end{bmatrix}$$

(iii) eigen vector of eigen value $\lambda = 2 - \sqrt{13}$

$$\lambda E - A = \begin{bmatrix} 1-\sqrt{13} & -3 & -6 \\ -2 & 1-\sqrt{13} & -4 \\ -1 & 0 & 1-\sqrt{13} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1-\sqrt{13} \\ 2 \\ 1 \end{bmatrix}$$

(5) eigen decomposition

$$A = \begin{bmatrix} 1 & 3 & 6 \\ 2 & 1 & 4 \\ 1 & 0 & 3 \end{bmatrix} \quad \Lambda = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2+\sqrt{13} & 0 \\ 0 & 0 & 2-\sqrt{13} \end{bmatrix} \quad Q = \begin{bmatrix} 2 & -1+\sqrt{13} & -1-\sqrt{13} \\ 2 & 2 & 2 \\ -1 & 1 & 1 \end{bmatrix}$$

$$\begin{vmatrix} 2 & 2 \\ 1 & 1 \end{vmatrix} = 0 \quad \begin{vmatrix} 2 & 2 \\ -1 & 1 \end{vmatrix} = 4 \quad \begin{vmatrix} 2 & 2 \\ -1 & 1 \end{vmatrix} = 4$$

$$\begin{vmatrix} -1+\sqrt{13} & -1-\sqrt{13} \\ 1 & 1 \end{vmatrix} = 2\sqrt{13} \quad \begin{vmatrix} 2 & -1-\sqrt{13} \\ -1 & 1 \end{vmatrix} = 1-\sqrt{13} \quad \begin{vmatrix} 2 & -1+\sqrt{13} \\ -1 & 1 \end{vmatrix} = 1+\sqrt{13}$$

$$\begin{vmatrix} -1+\sqrt{13} & -1-\sqrt{13} \\ 2 & 2 \end{vmatrix} = 4\sqrt{13} \quad \begin{vmatrix} 2 & -1-\sqrt{13} \\ 2 & 2 \end{vmatrix} = 2-2\sqrt{13} \quad \begin{vmatrix} 2 & -1+\sqrt{13} \\ 2 & 2 \end{vmatrix} = 6-2\sqrt{13}$$

$$Q^{-1} = \begin{bmatrix} 0 & 4 & 4 \\ \sqrt{13} & 1-\sqrt{13} & 1+\sqrt{13} \\ 4\sqrt{13} & 2-2\sqrt{13} & 6-2\sqrt{13} \end{bmatrix}$$

$$A = Q \Lambda Q^{-1}$$

$$\begin{bmatrix} 1 & 3 & 6 \\ 2 & 1 & 4 \\ 1 & 0 & 3 \end{bmatrix} = \begin{bmatrix} 2 & -1+\sqrt{13} & -1-\sqrt{13} \\ 2 & 2 & 2 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2+\sqrt{13} & 0 \\ 0 & 0 & 2-\sqrt{13} \end{bmatrix} \begin{bmatrix} 0 & 4 & 4 \\ \sqrt{13} & 1-\sqrt{13} & 1+\sqrt{13} \\ 4\sqrt{13} & 2-2\sqrt{13} & 6-2\sqrt{13} \end{bmatrix}$$

Question 3 (Pandas, 20 pts)

In this question, you will be using **Pandas** to apply exploratory data analysis of a Covid-19 dataset (from **The New York Times**).

If you have not installed the required packages, please refer to the **lab session material** for instructions.

Reading data using Pandas

```
In [ ]: # Read the dataset you will be working on
# The dataframe loaded with pandas is named as data
import pandas as pd
data = pd.read_csv('covid_19.csv')

# Take a look at the first 3 rows
data.head(3)
```

```
Out[ ]:
```

	date	state	cases	deaths
0	1/21/20	Washington	1	0
1	1/22/20	Washington	1	0
2	1/23/20	Washington	1	0

Question 3.1 (Get the shape of data, 5 pts)

Print the number of rows and columns of the dataframe "data"

```
In [ ]: ##### Your answer for Question 3.1 #####
##### Your code here #####
import pandas as pd
data = pd.read_csv('covid_19.csv')
print("number of rows:", data.shape[0])
print("number of columns:", data.shape[1])

number of rows: 31089
number of columns: 4
```

Data information

In Pandas, there are many summary functions which contain statistics as well as other data information. The name of the columns are:

```
In [ ]: data.columns

Out[ ]: Index(['date', 'state', 'cases', 'deaths'], dtype='object')
```

A brief summary of the dataset information:

mean: Mean of the values.

std: Standard deviation of the observations.

25%: The lower percentile.

75%: The upper percentile.

You may use `.describe()` to get a brief summary of the dataframe information.

```
In [ ]: data.describe()
```

```
Out[ ]:
```

	cases	deaths
count	3.108900e+04	31089.000000
mean	3.235684e+05	6171.822413
std	5.600332e+05	10224.348148
min	1.000000e+00	0.000000
25%	1.670600e+04	362.000000
50%	1.108810e+05	2075.000000
75%	4.098610e+05	7360.000000
max	4.647180e+06	68034.000000

To show the summarized information of a variable (i.e., the variable "deaths"):

```
In [ ]: # We can access a certain variable ("deaths") of the dataframe ('data') simply through data.deaths
data.deaths.describe()
```

```
Out[ ]: count    31089.000000
mean      6171.822413
std       10224.348148
min         0.000000
25%        362.000000
50%       2075.000000
75%       7360.000000
max       68034.000000
Name: deaths, dtype: float64
```

Missing values and data types

Entries with missing values are usually assigned with the value **NaN** ("Not a Number"), and the datatype is float64 dtype.

Question 3.2 (Check missing values, 5 pts)

Check whether there are missing values in the dataframe: print how many missing values exist in each column.

```
In [ ]: ##### Hint: this dataset does not have empty values #####
##### Your answer for Question 3.2 #####
##### Your code here #####
datenull = data.date.isnull()
statenull= data.state.isnull()
casesnull = data.cases.isnull()
deathsnull = data.deaths.isnull()
print("Number of missing values in column date:", datenull.sum())
print("Number of missing values in column state:", statenull.sum())
print("Number of missing values in column cases:", casesnull.sum())
print("Number of missing values in column deaths:", deathsnull.sum())

Number of missing values in column date: 0
Number of missing values in column state: 0
Number of missing values in column cases: 0
Number of missing values in column deaths: 0
```

Indexing and slicing

```
In [ ]: # Get the 10-th row for variable "State"
data['state'][10]
```

```
Out[ ]: 'Illinois'
```

Index based selection with `iloc`: `iloc` is **row-first, column-second**.

```
In [ ]: # The first row of the dataframe
print(data.iloc[0])
```

```

date          1/21/20
state         Washington
cases          1
deaths         0
Name: 0, dtype: object

```

```

In [ ]: # The first column of the dataframe
print(data.iloc[:, 0])

0      1/21/20
1      1/22/20
2      1/23/20
3      1/24/20
4      1/24/20
...
31084   9/18/21
31085   9/18/21
31086   9/18/21
31087   9/18/21
31088   9/18/21
Name: date, Length: 31089, dtype: object

```

```

In [ ]: # The first column (and 2nd-5th rows) of the dataframe
print(data.iloc[2:6, 0])
# or pass a list
print(data.iloc[[i+2 for i in range(4)], 0])

2      1/23/20
3      1/24/20
4      1/24/20
5      1/25/20
Name: date, dtype: object
2      1/23/20
3      1/24/20
4      1/24/20
5      1/25/20
Name: date, dtype: object

```

Question 3.3 (Conditional selection, 5 pts)

What are the number of "cases" and "deaths" for 'California' on '8/21/21'? (print the corresponding row in this dataframe with `loc`)

```

In [ ]: ##### Your answer for Question 3.3 #####
#data.loc[]
##### Your code here (complete the code above) #####
newdata = data.loc[data["state"] == "California"][data["date"] == "8/21/21"]
print(newdata.cases.sum(), newdata.deaths.sum())

```

```
4316350 65082
```

C:\Users\rinbe\AppData\Local\Temp\ipykernel_22844\4164158780.py:4: UserWarning: Boolean Series key will be reindexed to match DataFrame index.

```
newdata = data.loc[data["state"] == "California"][data["date"] == "8/21/21"]
```

Question 3.4 (Data aggregation, 5 pts)

Add a new column named "ratio" (for the dataframe "data") which defined as the ratio "deaths"/"cases" in each row.

```

In [ ]: ##### Your answer for Question 3.4 #####
##### Your code here #####
data["ratio"] = data["deaths"]/data["cases"]
data.tail()

```

```
Out[ ]:
```

	date	state	cases	deaths	ratio
31084	9/18/21	Virginia	827197	12242	0.014799
31085	9/18/21	Washington	623254	7256	0.011642
31086	9/18/21	West Virginia	221513	3370	0.015214
31087	9/18/21	Wisconsin	772089	8703	0.011272
31088	9/18/21	Wyoming	83958	918	0.010934

Question 4 (Seaborn and Matplotlib, 30 pts)

Visualizing pairplots using seaborn

Seaborn: Python library for statistical data visualization built on top of Matplotlib

Tutorial: detailed example codes are [here](#) if needed.

Now we shortly switch our focus to data that only about California, Arizona and Washington.

```
In [ ]: # the sub-dataframe contains only 'California', 'Arizona', 'Washington' is named as subset
subset = data.loc[data['state'].isin(['California', 'Arizona', 'Washington'])]
subset = subset.reset_index(drop=True)
subset.head()
```

```
Out[ ]:
```

	date	state	cases	deaths	ratio
0	1/21/20	Washington	1	0	0.0
1	1/22/20	Washington	1	0	0.0
2	1/23/20	Washington	1	0	0.0
3	1/24/20	Washington	1	0	0.0
4	1/25/20	California	1	0	0.0

```
In [ ]: # import required packages
import seaborn as sns
import matplotlib.pyplot as plt
# Allow figures to be shown in the jupyter notebook interface
%matplotlib inline
```

Question 4.1 (Visualizing statistical relationships, 10 pts)

In Seaborn, `relplot()` provides access to several different axes-level functions that show the relationship between two variables with semantic mappings of subsets.

Adopt `relplot()` and visualize how the variable "cases" changes w.r.t "date" for three selected states.

A basic tutorial is [here](#).

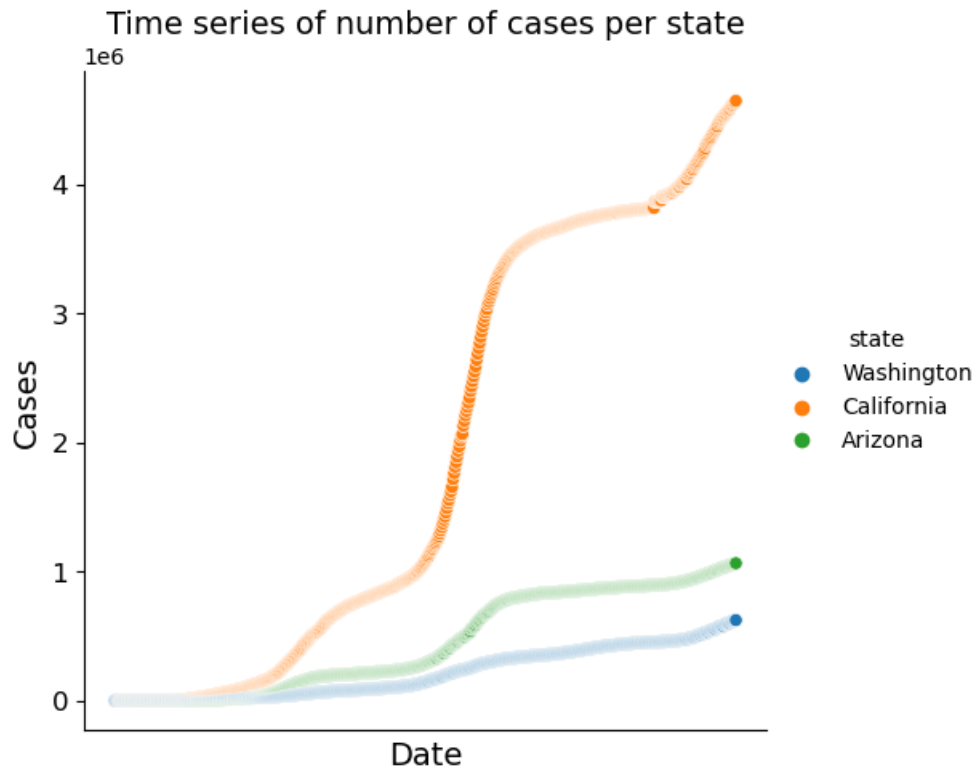
```
In [ ]: ##### Your answer for Question 4.1 #####

##### Your code here #####
sns.relplot(data=subset, x=subset["date"], y=subset["cases"], hue=subset['state'])

# Ignore the xticks (since too many dates)
plt.tick_params(
    axis='both',          # changes apply to the x-axis
    which='both',        # both major and minor ticks are affected
    bottom=False,         # ticks along the bottom edge are off
    top=False,            # ticks along the top edge are off
    labelbottom=False
)
```

```
# Use matplotlib to modify figure parameters
plt.ylabel('Cases', fontsize=14)
plt.yticks(fontsize=12)
plt.xlabel('Date', fontsize=14)
plt.title('Time series of number of cases per state', fontsize=14)
plt.show()
```

c:\Users\rinbe\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout has changed to tight
self._figure.tight_layout(*args, **kwargs)



Question 4.2 (Regression plot with Seaborn, 5 * 4 pts)

In Seaborn, there are several statistical models to estimate a simple relationship between two sets of observations. A basic tutorial is [here](#).

Suppose we are only interested in covid-19 information for "California"

```
In [ ]: # Only adopt samples of California information

data_ca = subset[subset['state']=='California']
print(data_ca.shape)
data_ca = data_ca.reset_index(drop=True)
data_ca.head()
```

(603, 4)

```
Out[ ]:   date    state  cases  deaths
0  1/25/20  California    1      0
1  1/26/20  California    2      0
2  1/27/20  California    2      0
3  1/28/20  California    2      0
4  1/29/20  California    2      0
```

Question 4.2.1 Visualize 1 (5 pts)

Use seaborn `regplot()` to visualize the relationship between "date_order" "deaths". (select only first 50 rows of the dataframe "data_ca")

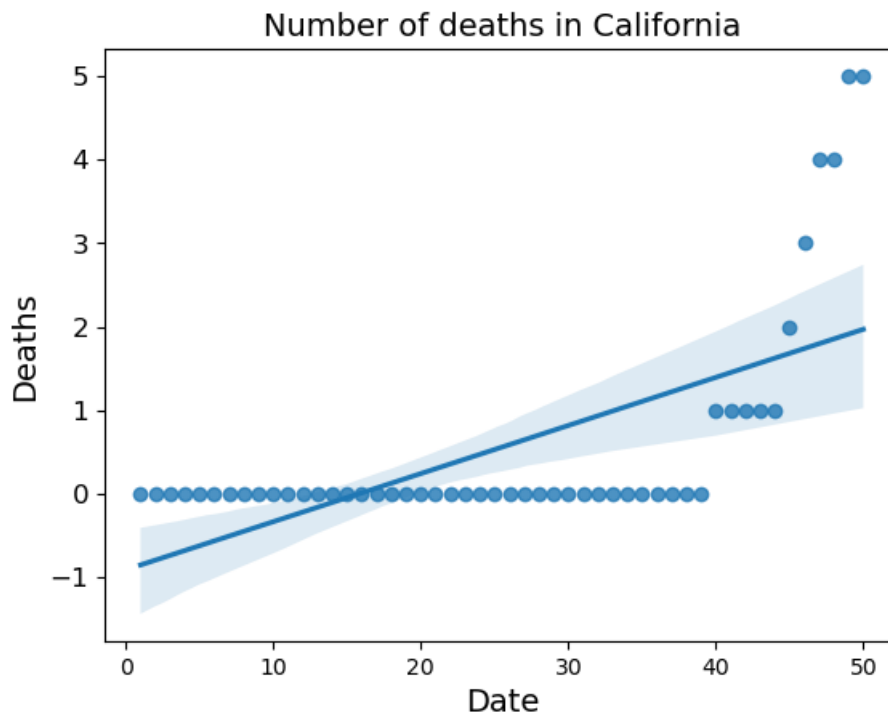
A reference link is [here](#).

```
In [ ]: data_ca['date_order'] = [i+1 for i in range(data_ca.shape[0])]

##### Your answer for Question 4.2.1 #####
dataca2 = data_ca.iloc[0:50]
sns.regplot(data=dataca2,x=dataca2["date_order"], y=dataca2["deaths"], x_estimator=None, x_bins=None, x_ci='ci')

##### Your code here (reminder: select "date_order" rather than "date") #####

plt.ylabel('Deaths', fontsize=14)
plt.yticks(fontsize=12)
plt.xlabel('Date', fontsize=14)
plt.title('Number of deaths in California', fontsize=14)
plt.show()
```



Question 4.2.2 Visualize 2 (5 pts)

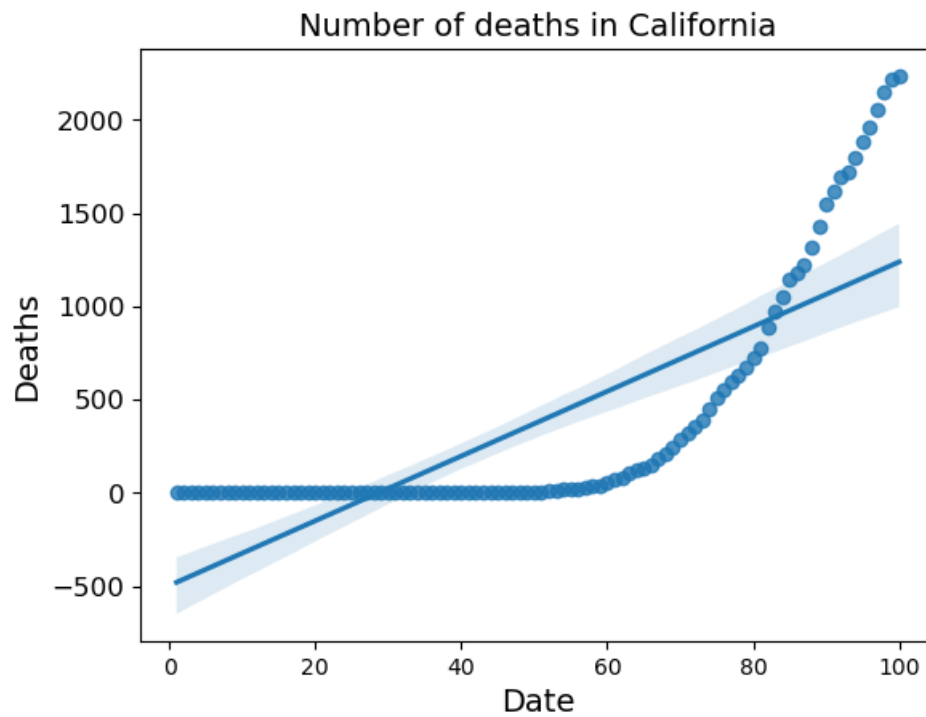
Use seaborn `regplot()` to visualize the relationship between "date_order" "deaths". (select first 100 rows of the dataframe "data_ca")

Same as Visualize 1, but with more rows of the dataframe "data_ca" included. A reference link is [here](#).

```
In [ ]: ##### Your answer for Question 4.2.2 #####

##### Your code here (only visualize w.r.t. first 100 rows of dataframe "data_ca") #####
dataca3 = data_ca.iloc[0:100]
sns.regplot(data=dataca3,x=dataca3["date_order"], y=dataca3["deaths"], x_estimator=None, x_bins=None, x_ci='ci')

plt.ylabel('Deaths', fontsize=14)
plt.yticks(fontsize=12)
plt.xlabel('Date', fontsize=14)
plt.title('Number of deaths in California', fontsize=14)
plt.show()
```



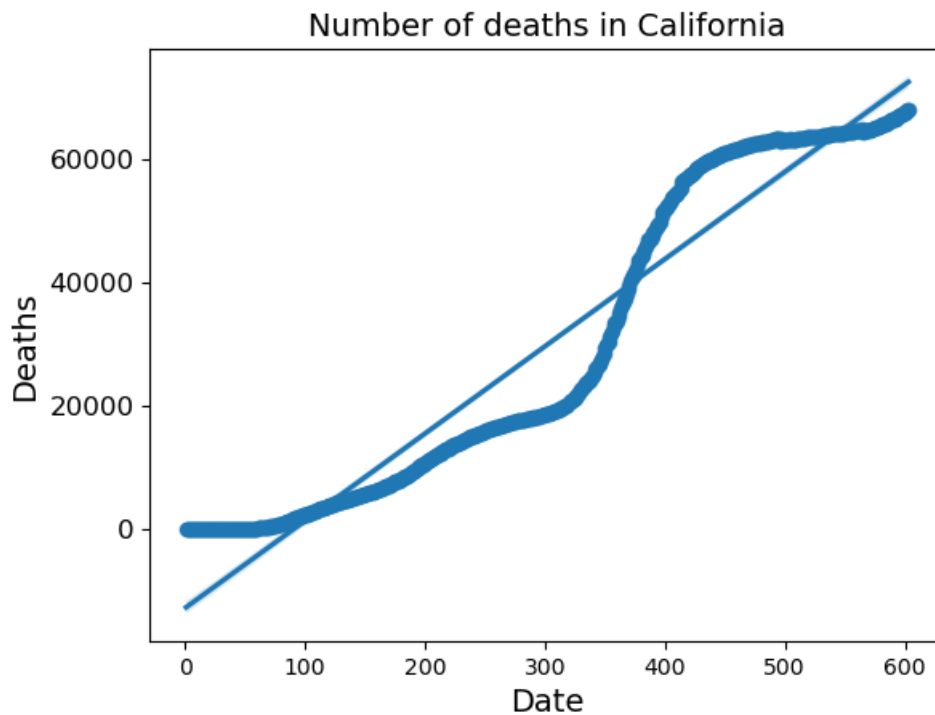
Question 4.2.3 Visualize 3 (5 pts)

Use seaborn `regplot()` to visualize the relationship between "date_order" "deaths". (use the whole dataframe "data_ca")

Same as Visualize 1, but use the whole dataframe "data_ca".

```
In [ ]: ##### Your answer for Question 4.2.3 #####

##### Your code here #####
sns.regplot(data=data_ca, x=data_ca["date_order"], y=data_ca["deaths"], x_estimator=None, x_bins=None, x_ci='ci',
plt.ylabel('Deaths', fontsize=14)
plt.yticks(fontsize=12)
plt.xlabel('Date', fontsize=14)
plt.title('Number of deaths in California', fontsize=14)
plt.show()
```

Question 4.2.4 What is your observations from the above three figures? (5 pts, open question)

- The observed dots on the visualized plot became non-contiguous to continuous as the number of data frames (or population) increased. As the dots become continuous, the observed points become closer to the estimated linear curve as well.
- There had been a visually dynamic increase in the plot for 50 rows, however, this was a non-dynamic continuous plot when seen in the long term. From the plot for the whole data frame, the data makes an increase around day 100 for the first time and gradually increases with the estimate line.

From the observation, it could be said that the number of data frames affects the accuracy of the estimated plot. The greater the number of data, the more accurate the estimated plot will remain. However, when using a linear curve for the estimated plot, this is unlikely to be accurate than other multi-vector plots. This can be seen from the fact that an event exists though the number of deaths is not likely to fall below zero.