



EPA Source Apportionment Fall 2025 Capstone

Deeya Datta, Prateek Mishra, & George Whittington
Client: Deron Smith, EPA



Table of contents

O1

Introduction

Background information and client

O2

Data

The data given for this project and questions

O3

EDA

Exploratory data analysis

O4

Future Work

Next steps for next semester





01

Introduction

Our Project Topic and Client





Source Apportionment

Source apportionment is a "collection of techniques to provide information regarding how much a source (usually a generalized category) contributes to the overall pollutant concentration at receptor (usually a monitoring site) (Rizzo, 2010).

The overall goal is to extract latent factors from the data. We can then interpret these factors to identify specific source profiles.





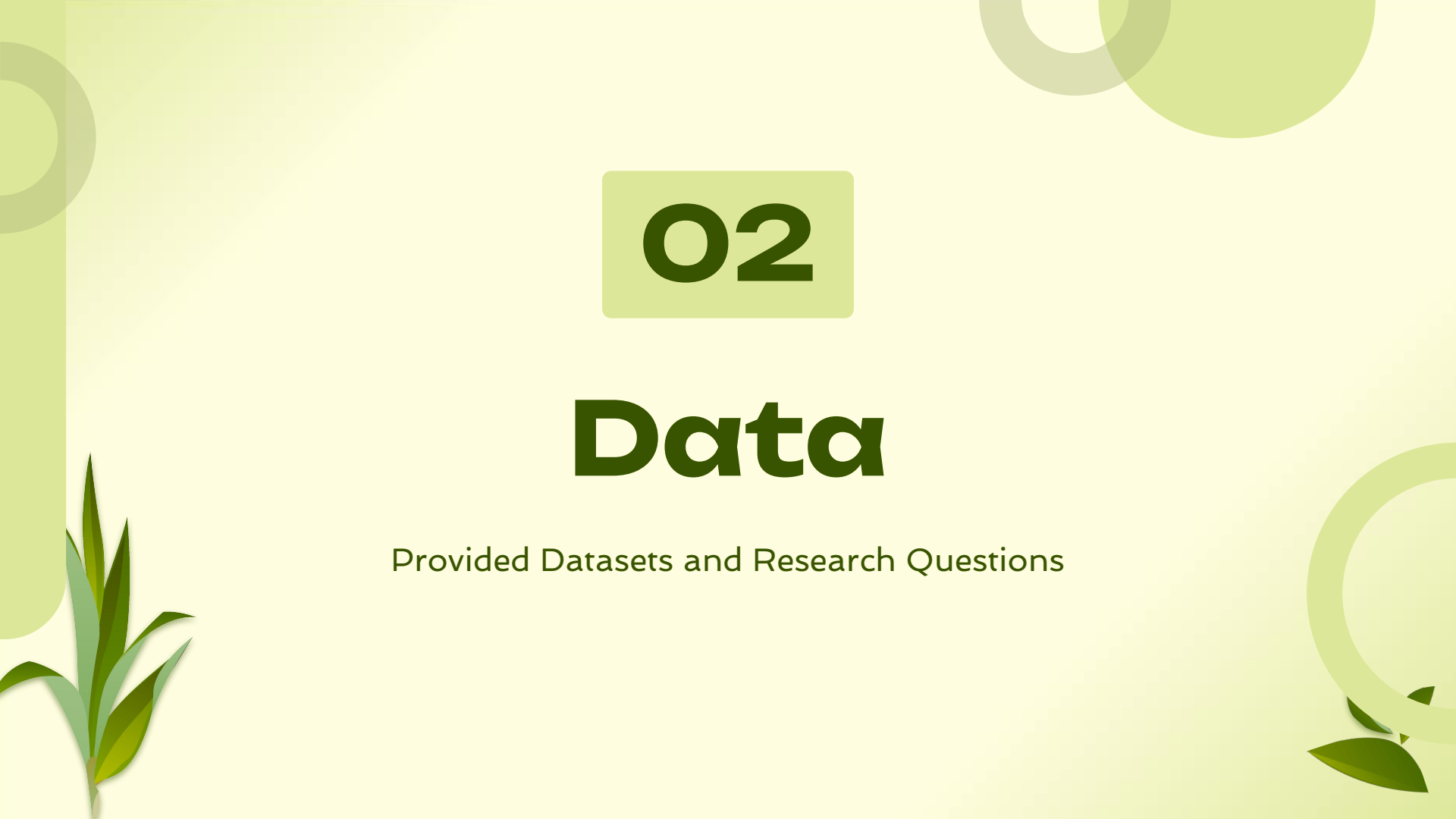
Deron Smith



Mr. Smith is a UGA alumni with a B.S. in Physics and a M.S. in Computer Science from GaTech. He is now a researcher and software developer for the Environmental Protection Agency

Recent projects include an open source tool called the Environmental Source Apportionment Toolkit (esat), and is currently working on a graphical user interface for the toolkit.



The slide features a light green background with decorative elements. In the top left, there is a large, faint green circle. In the top right, there are two overlapping green circles. In the bottom left, there is a green leafy plant. In the bottom right, there is a green leafy plant and a large, faint green circle. The number '02' is displayed in a bold, dark green font inside a light green rounded square.

02

Data

Provided Datasets and Research Questions



The Datasets



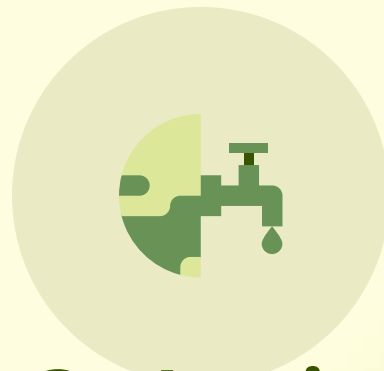
Baltimore

focuses on fine
particulate matter
(PM_{2.5}) chemical
composition (n=630)



Baton Rouge

focuses on volatile
organic compounds
(VOCs) rather than
particulate matter
(n=307)



St. Louis

focuses on fine
particulate matter
(PM_{2.5}) chemical
composition (n=418)





Variables Across Cities

City	Category	Count	Specific Species Included
Baltimore	Trace Elements	17	Al, As, Ba, Br, Ca, Cl, Cr, Cu, Fe, Mn, Ni, Pb, Se, Si, Ti, V, Zn
	Ions	5	NH_4^+ , K^+ , Na^+ , SO_4^{2-} , NO_3^-
	Carbon & Mass	6	EC, OC, OM, PM2.5
Baton Rouge	VOCs	41	BTEX, Alkanes, Alkenes & Isoprene
St. Louis	Trace Elements	8	Cd, Cu, Fe, Mn, Ni, Pb, Se, Zn
	Ions	2	SO_4^{2-} , NO_3^-
	Carbon & Mass	3	EC, OC, PM2.5





Considerations & Challenges

Temporal Resolution & Irregularity

Baltimore:

- Irregular daily measurements with frequent multi-day gaps.

Baton Rouge:

- Sparse hourly sampling (mostly 03:00 or 06:00) rather than continuous monitoring.

St. Louis:

- "Chunked" hourly data (e.g., 00:00–23:00) but with several intraday gaps (missing afternoon blocks).





Considerations & Challenges

Completeness & Detection Limits

Data Integrity:

- There were no missing values found across any of the three datasets.

Signal-to-Noise:

- However, we identified observations where the instrument uncertainty exceeded the measured concentration.

Handling Low Values:

- These "high-noise" points are treated as Below Detection Limit (BDL) but preserved using log-transformation.





Research Questions

Input Data Analysis

How can the input data be effectively analyzed and pre-processed to identify key features?

Source Apportionment Modeling

Can we calculate latent Source Profiles and quantify their Contributions using a standard apportionment model?

Uncertainty Integration

How can we utilize the paired measurement uncertainty into the model to refine our estimates and assess potential error?

Model Evaluation

How can we evaluate the results of multiple models (metrics, plots) to ensure stability?

Environmental Interpretation

What physical conclusions can we make about the real-world pollution sources based on the model outputs?

The slide features a light green background with decorative elements. In the top left, there is a large, faint green circle. In the top right, there are two overlapping green circles. In the bottom left, there is a green leafy plant. In the bottom right, there is a green leafy plant and a large, faint green circle.

03

EDA

Exploratory Data Analysis



Gaps Between Extreme Values

Species	mean	max	max_mean_ratio
PM2.5	15.57	76.30	4.90
OM	6.60	58.10	8.80
Organic Carbon	4.71	41.50	8.80
Sulfate	4.83	30.20	6.25
Total Nitrate	1.83	12.60	6.88

Baltimore

Baton Rouge

Species	mean	max	max_mean_ratio
TNMOC	199.90	708.47	3.54
Unidentified	63.20	385.18	6.09
Isopentane	16.40	125.50	7.65
N-Hexane	8.54	92.56	10.84
Propane	15.78	70.96	4.50

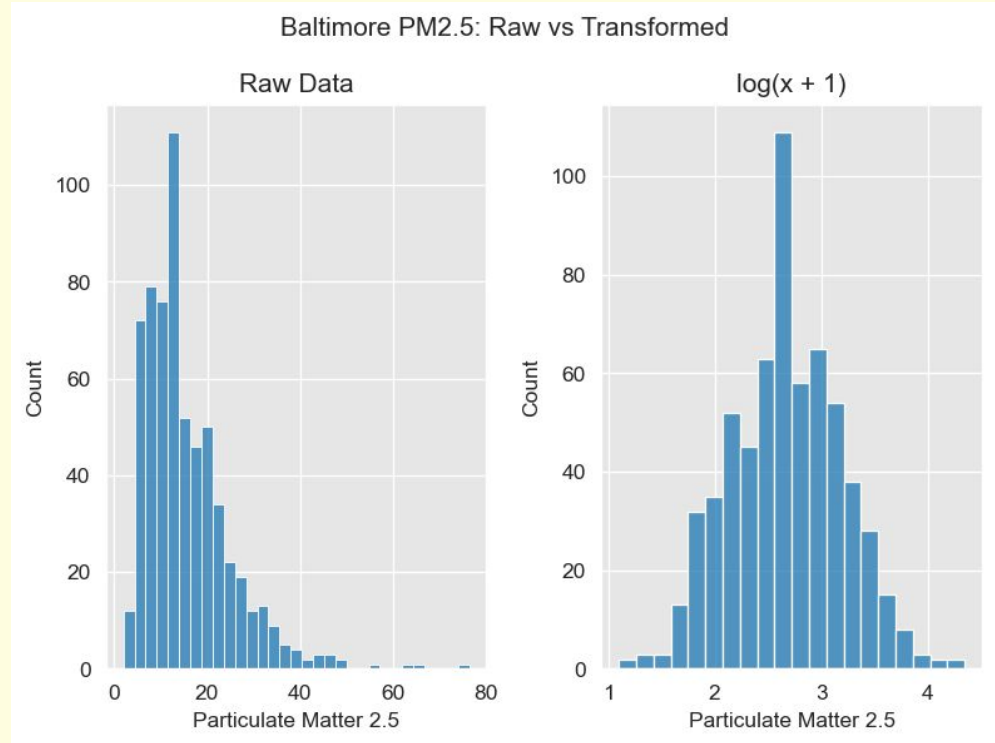
Species	mean	max	max_mean_ratio
Mass	22.95	107.00	4.66
OC	4.77	71.20	14.91
SO4	3.84	16.90	4.40
NO3	1.82	10.30	5.65
EC	0.77	6.38	8.28

St. Louis



$\log(x + 1)$ Transformation

A $\log(x + 1)$ is applied to the data to try and bring the data to be more symmetric and to also preserve the non-negative status

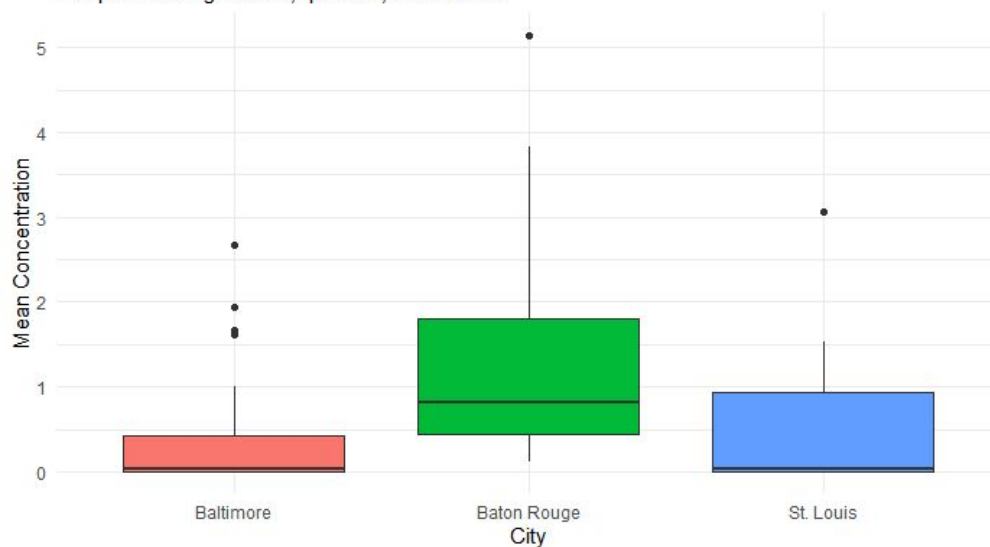




Concentration Distributions

Distribution of Mean Concentrations

Box plot showing median, quartiles, and outliers



TOP 5 SPECIES BY MEAN CONCENTRATION

BATON ROUGE (VOCs in ppbc):

1. <code>tnmoc</code>	5.1488
2. <code>unidentified</code>	3.8296
3. <code>propane</code>	2.6409
4. <code>ethane</code>	2.5839
5. <code>isopentane</code>	2.5333

ST. LOUIS ($\mu\text{g}/\text{m}^3$):

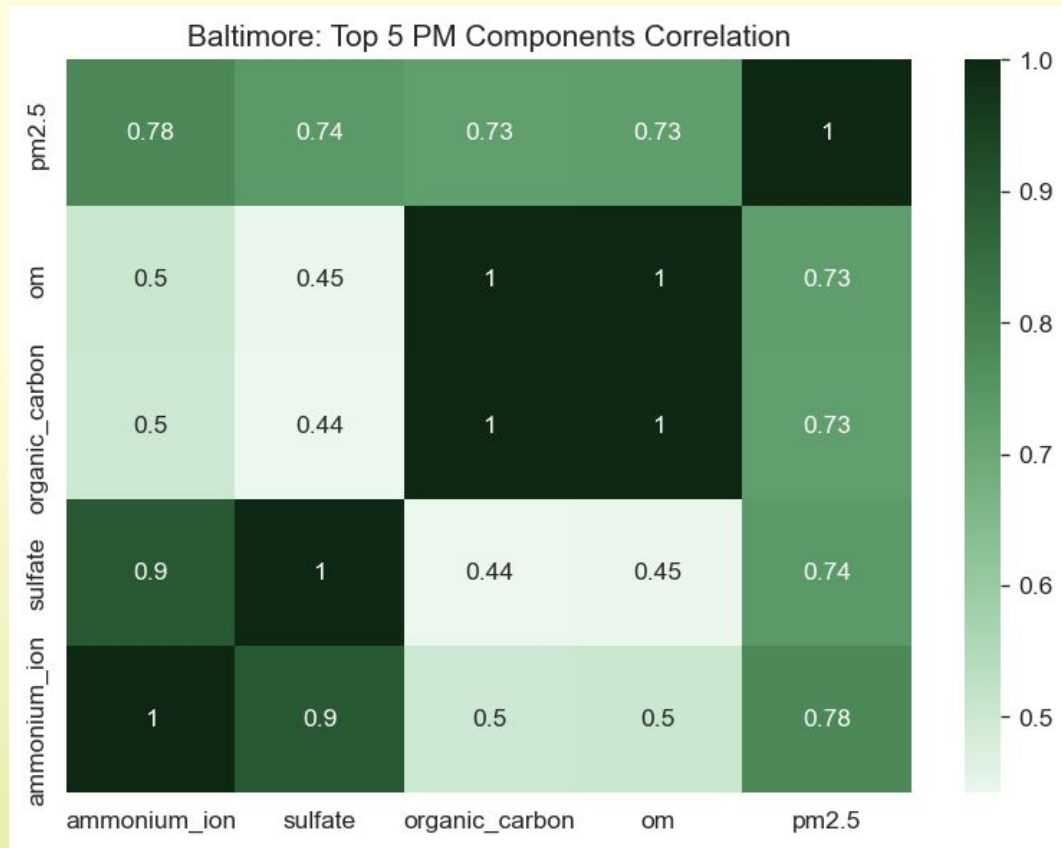
1. <code>pm2.5</code>	3.0585
2. <code>organic_carbon</code>	1.5261
3. <code>sulfate</code>	1.4111
4. <code>total_nitrate</code>	0.9303
5. <code>elemental_carbon</code>	0.5190

BALTIMORE ($\mu\text{g}/\text{m}^3$):

1. <code>pm2.5</code>	2.6672
2. <code>om</code>	1.9389
3. <code>organic_carbon</code>	1.6625
4. <code>sulfate</code>	1.6125
5. <code>ammonium_ion</code>	1.0076



Top Correlations



Organic Carbon and Organic Matter have a perfect correlation of 1.0. This is because OM isn't independently measured, it's simply OC multiplied by a conversion factor to account for the hydrogen and oxygen attached to the carbon.



The slide features a light green background with decorative elements. In the top-left corner, there is a large, faint green circle. In the top-right corner, there are two overlapping green circles. In the bottom-left corner, there is a green leafy plant. In the bottom-right corner, there is a green leafy plant and a large, faint green circle.

O4

Future Work

Goals For Next Semester

How do we plan to implement a model?

What is Positive Matrix Factorization?

PMF separates pollutants into sources and their contributions.

Source Contributions (G): how much each source contributed on each day

Source Profiles (F): the chemical fingerprint of each source

$X = GF + E$ where X is observed and E are residual errors

Why is it a good fit?

- Our data is a mixture of many unknown sources

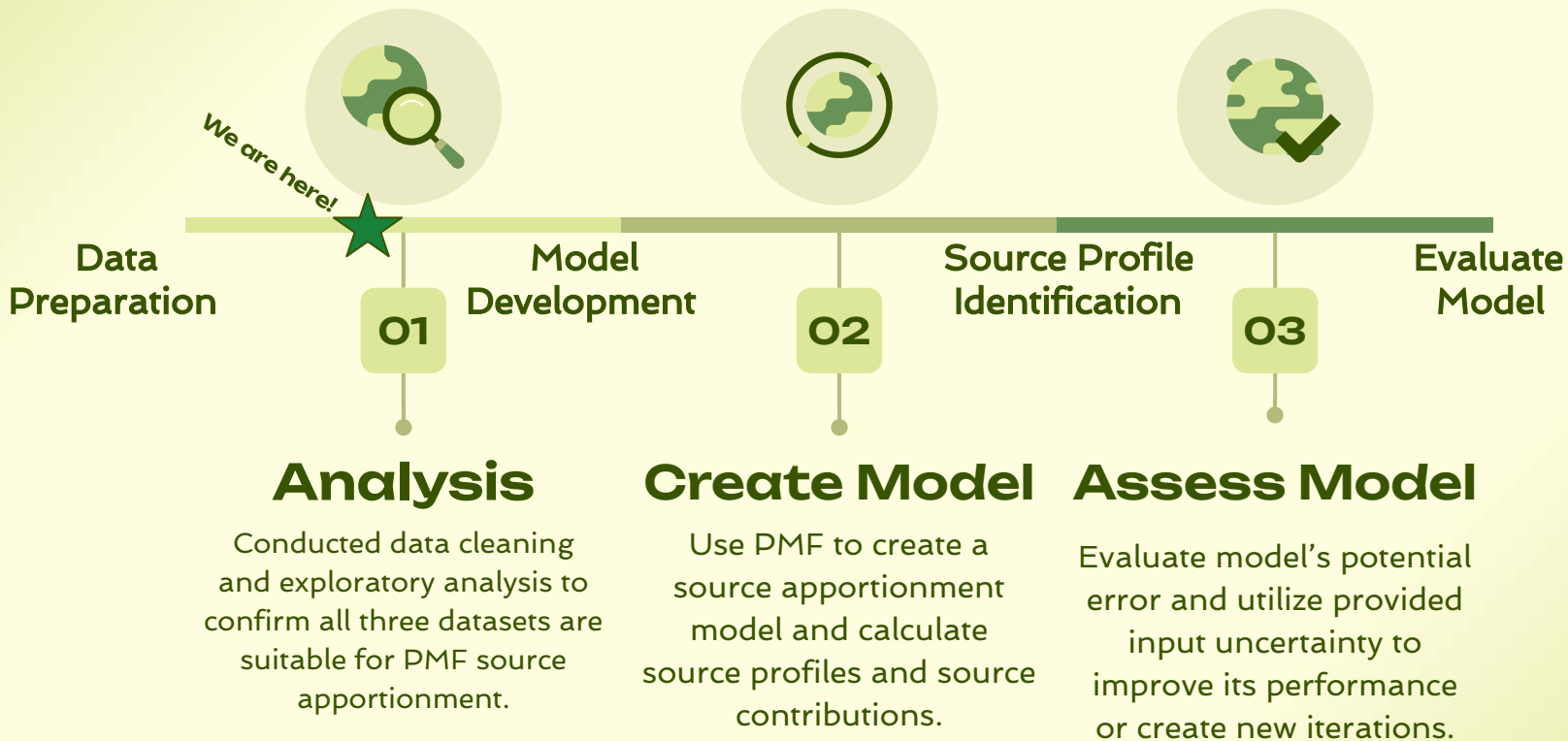
- PMF handles uncertainty values (low uncertainty are weighted higher)

- It requires no predefined source profiles

- Air quality researchers widely use it, ensuring reliability



Our Next Steps





Thanks!

Do you have any questions?

CREDITS: This presentation template was created by [Slidesgo](#), and includes icons by [Flaticon](#), and infographics & images by [Freepik](#)

