



IPS

Instituto
Politécnico de Setúbal
Escola Superior de
Tecnologia de Setúbal



Course Name:

Supervised Machine Learning

By Prof. Vala A. Rohani
vala.ali.rohani@estsetubal.ips.pt



Session outline

1. Introducing myself
2. Course Topics
3. Introduction to Data Analytics
4. Installing R and R Studio

Introducing myself

- Post doctorate in Data Science
- PhD in Software Engineering (Recommender Systems)
- University Lecturer for more than 10 years
- Having more than 20 published papers in index conferences and journals

Some of my Professional Certificates:

- Hadoop Platform and Application Framework from University of California
- Practical Machine Learning from Johns Hopkins University
- R Programming from Johns Hopkins University
- MongoDB for DBAs from MongoDB
- Social Network Analysis from University of Michigan
- Mining Massive Datasets from Stanford University
- Pattern Discovery in Data Mining from Illinois University
- Process Mining from Eindhoven University of Technology

**Vala Ali Rohani**

Invited Assistant Professor
Founder and Chief Data
Scientist at Data Corner
vala@data-corner.com



Course Topics



Course Topics

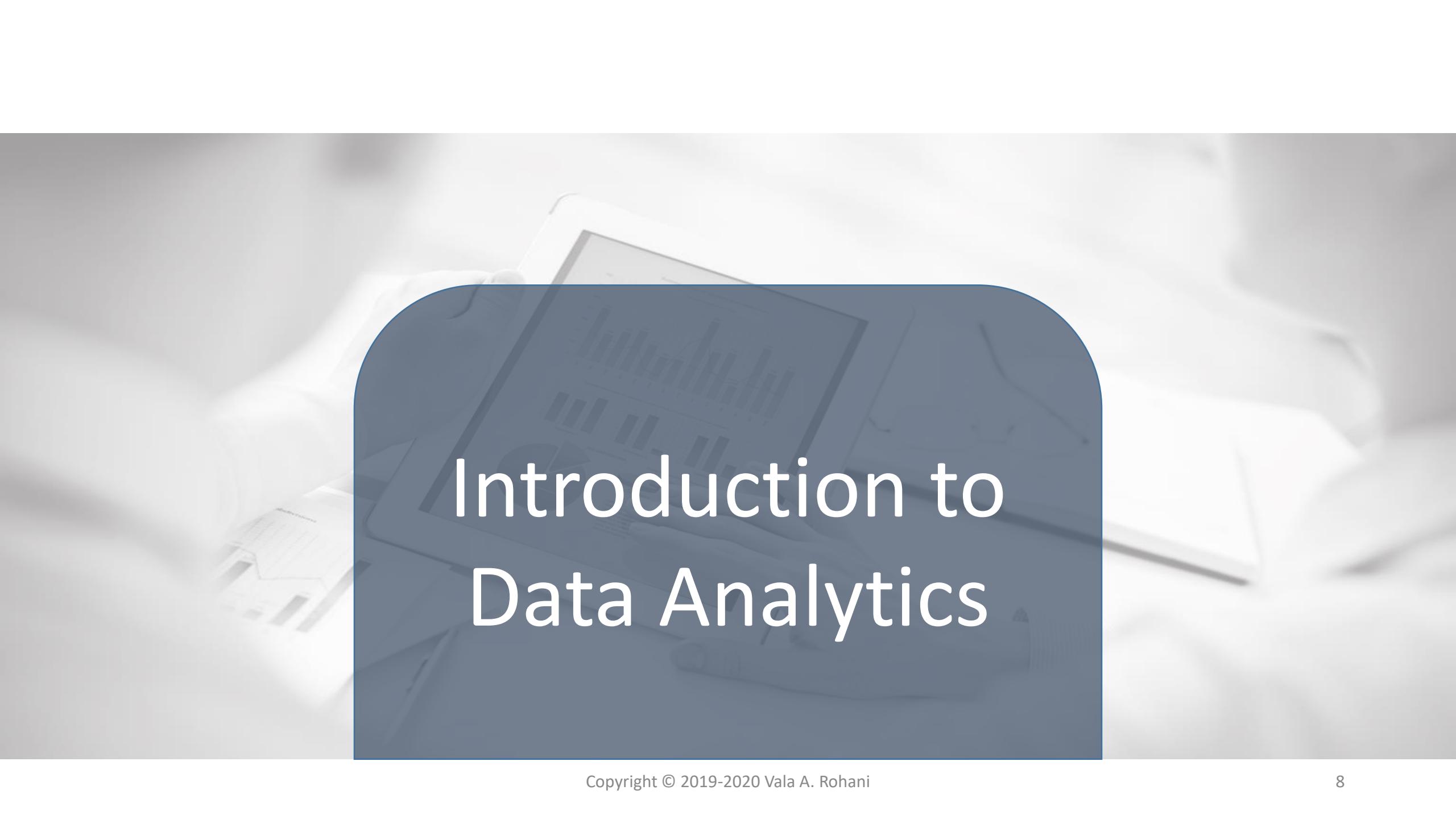
Docente	Sen	2ºF	3ºF	4ºF	5ºF	6ºF	Remoto/	Milestone			
VR	1 (4h)	08-Jan	09-Jan	10-Jan	11-Jan	12-Jan	P	Lab1	Presentation of course syllabus and evaluation schema. Introduction to Data Science and Machine Learning Reviewing the essential topics in R Programming Exploratory Data Analysis in R	Lab1 - Setting up the R and R Studio Reviewing the essential topics in R Programming Exploratory Data Analysis in R	
VR	2 (8h)	15-Jan	16-Jan	17-Jan	18-Jan	19-Jan	R	Lab2	Preparation and Data Cleaning for ML (Machine Learning) in R Preparing Training & Test datasets Regression Linear Models (1)	Lab2 - Data Cleaning in R Preparing Training & Test datasets Building Regression linear Models in R (1)	
VR	3 (12h)	22-Jan	23-Jan	24-Jan	25-Jan	26-Jan	R	Lab3	Regression linear Models (2) Evaluating the Performance of Regression Models	Lab3 - Building Regression linear Models in R (2) Applications	
VR	4 (16h)	29-Jan	30-Jan	31-Jan	01-Feb	02-Feb	R	Lab4	Classification Models (1) Logistic Regression	Lab4 - Building Logistic Regression Applications	
VR	5 (20h)	05-Feb	06-Feb	07-Feb	08-Feb	09-Feb	P	Lab5	Classification Models (2) Linear Discriminant Analysis	Lab5 - Building Classification Models in R (2) Applications	
HP	6 (24h)			07-Feb			R	Lab6	Classification Models (3) Quadratic Discriminant Analysis Naive Bayes Comparison of Classification Methods	Lab6 - Building Classification Models in R (3) Applications	
HP	7 (28h)	12-Feb	13-Feb	14-Feb	15-Feb	16-Feb	R	Lab7	Tree-Based Methods (1) Regression Trees	Lab7 -Building Tree-Based Methods (1) Applications	
HP	8 (32h)			14-Feb			P	Lab8	Tree-Based Methods (2) Classification Trees	Lab8 -Building Tree-Based Methods (2) Applications	
HP	9 (36h)	19-Feb	20-Feb	21-Feb	22-Feb	23-Feb	R	Lab9	Tree-Based Methods (3) Bagging, Random Forests, Boosting, and Bayesian Additive Regression Trees	Lab9 -Building Tree-Based Methods (3) Applications	
HP	10 (40h)			21-Feb			R	Lab10	Final Project Presentations		
HP	11 (44h)	26-Feb	27-Feb	28-Feb	29/fev	01-Mar	P	Test	Written Test		
VR		04-Mar	05-Mar	06-Mar	07-Mar	08-Mar	P	Exam	Written Exam		

Marking Scheme



Marking Scheme

Marking (Continues)	Labs	20%	
	project (grp)	30%	
	Written test	50%	min 9
Exams (I, II [III])	Practice	50%	
	Written	50%	min 9

A grayscale photograph of a person's hand holding a smartphone. The phone's screen displays several data visualizations, including bar charts and line graphs, suggesting a theme of data analysis or technology. The background is blurred.

Introduction to Data Analytics



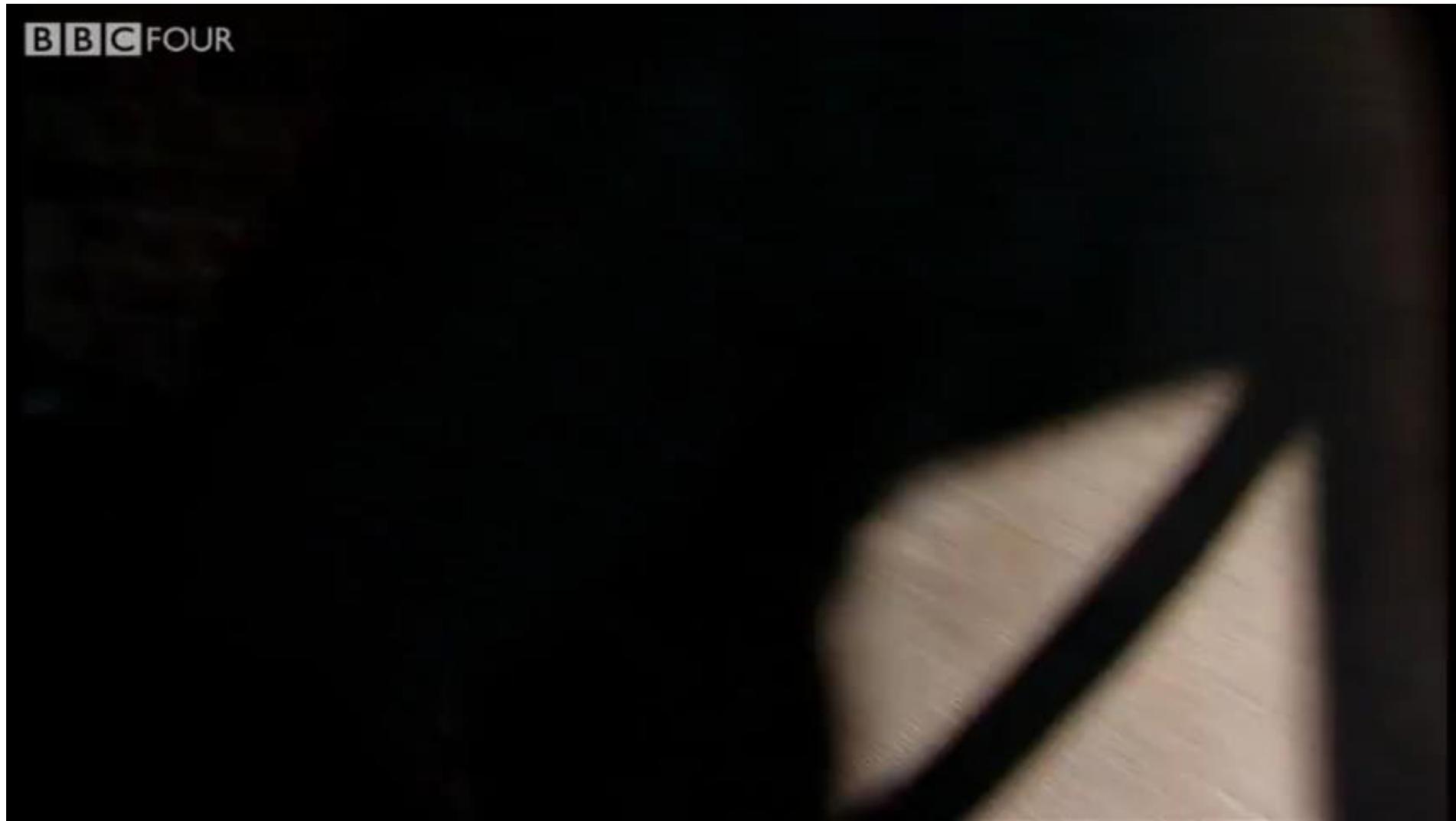
**INFORMATION IS THE OIL AND
ANALYTICS IS THE COMBUSTION ENGINE**

What is Data Analytics?

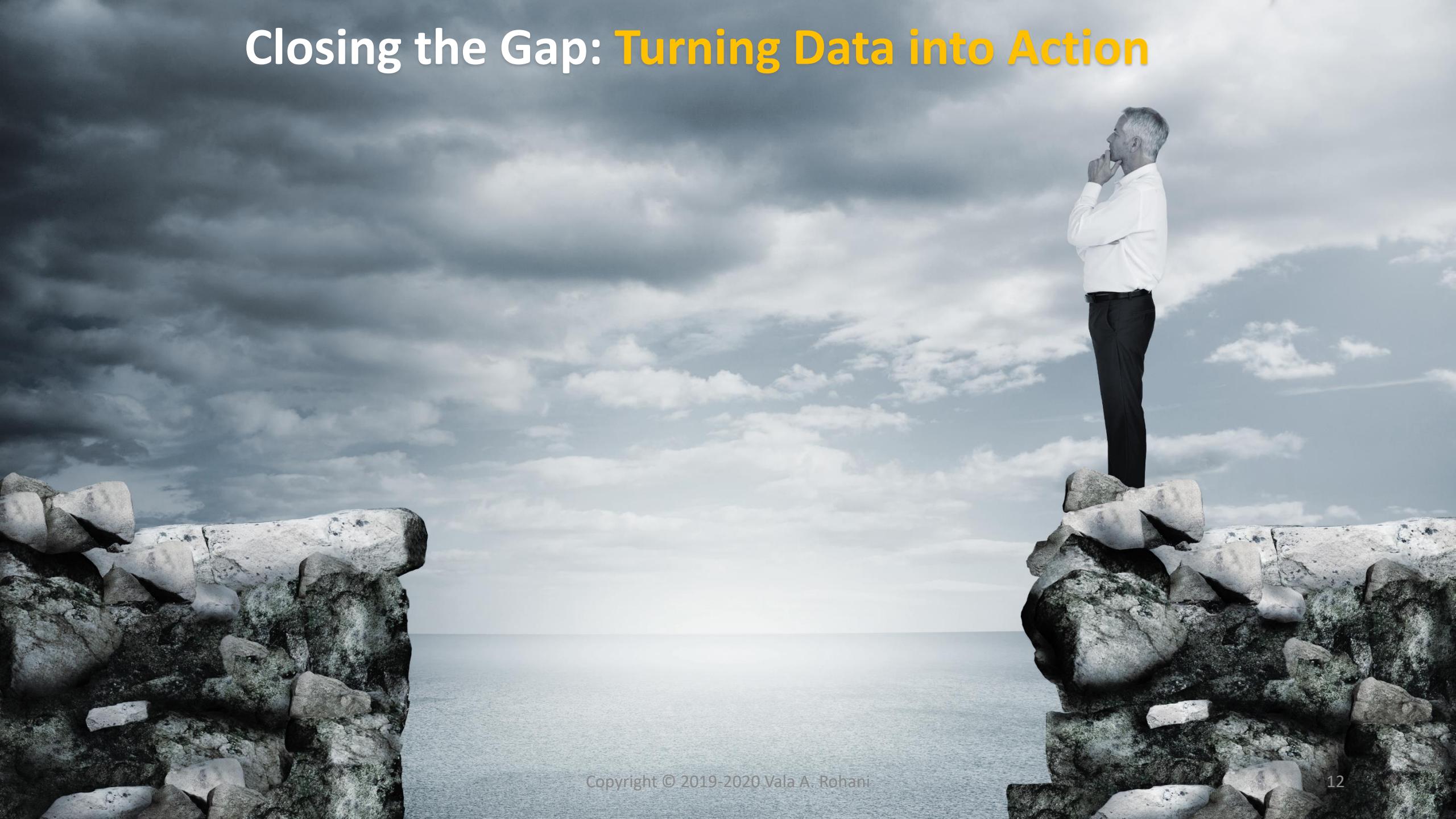
Data Analytics is the **discovery**, **interpretation**, and **communication** of meaningful patterns in data.

Analytics relies on the simultaneous application of **statistics**, **computer programming** and **operations research** to quantify performance and often favors data visualization to communicate insight.

Let's watch an inspired video on data visualization



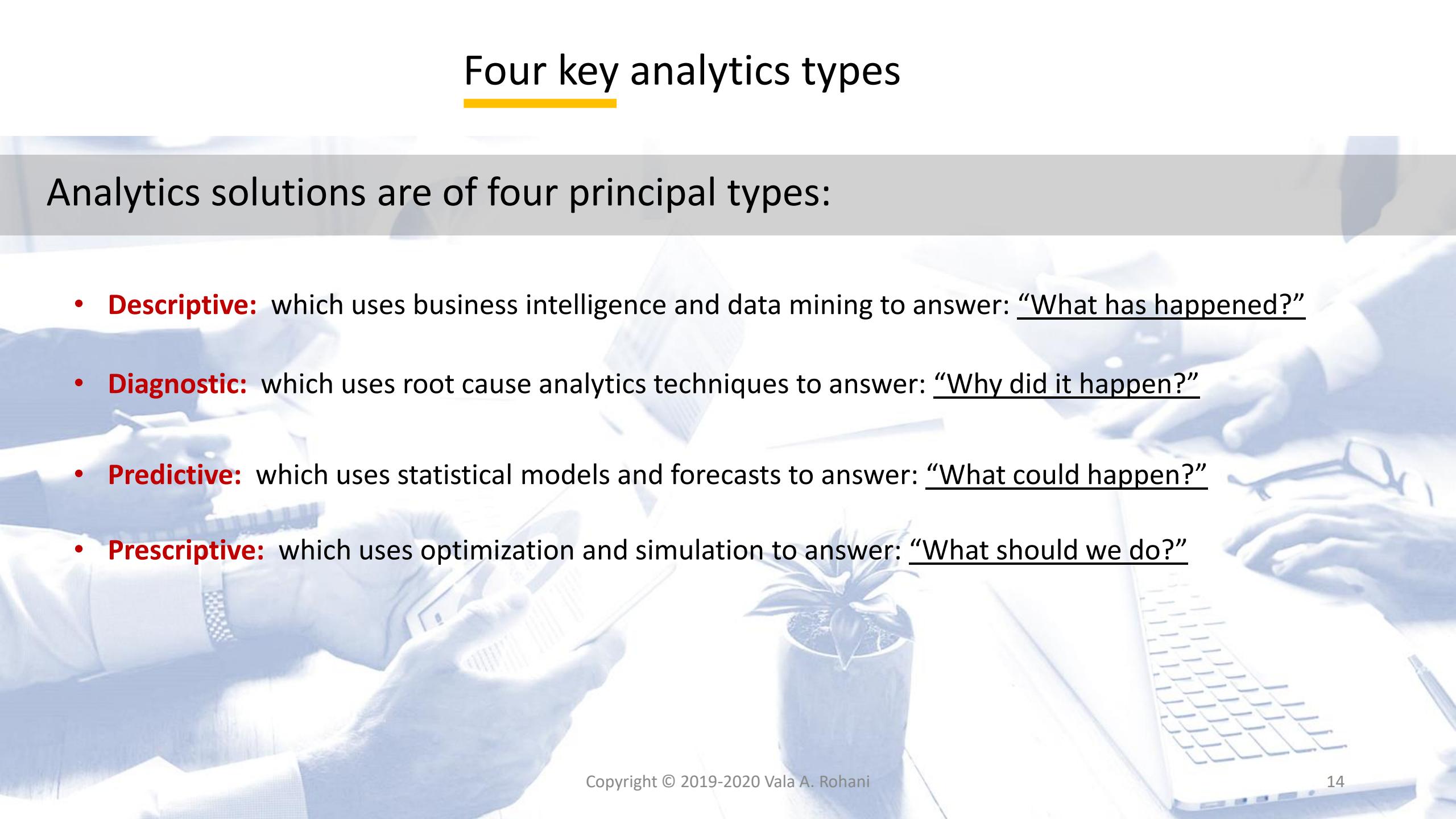
Closing the Gap: Turning Data into Action



Rise over the noise ...



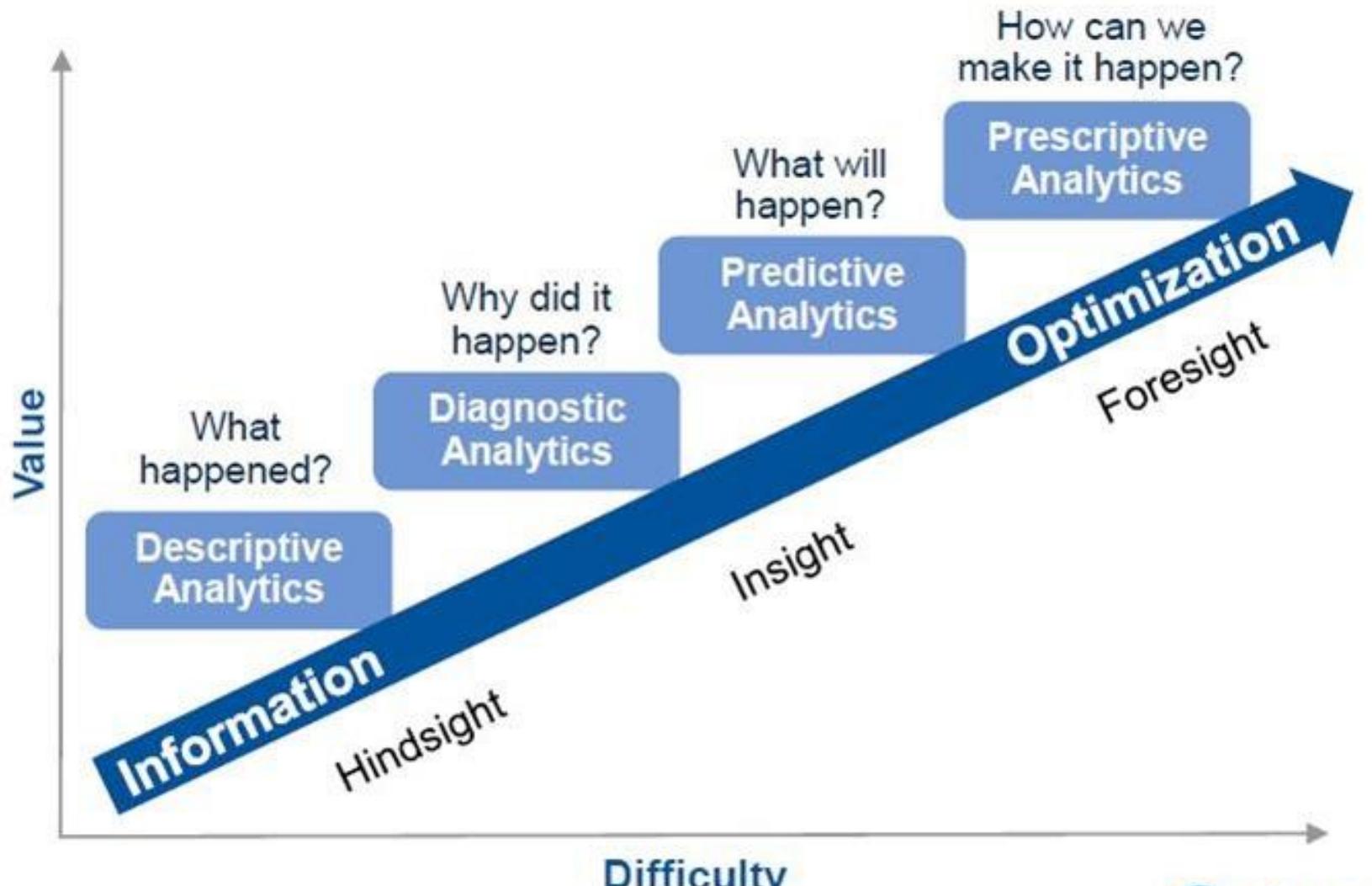
Four key analytics types



Analytics solutions are of four principal types:

- **Descriptive:** which uses business intelligence and data mining to answer: “What has happened?”
- **Diagnostic:** which uses root cause analytics techniques to answer: “Why did it happen?”
- **Predictive:** which uses statistical models and forecasts to answer: “What could happen?”
- **Prescriptive:** which uses optimization and simulation to answer: “What should we do?”

The analytics maturity model – Where are we in the journey ?

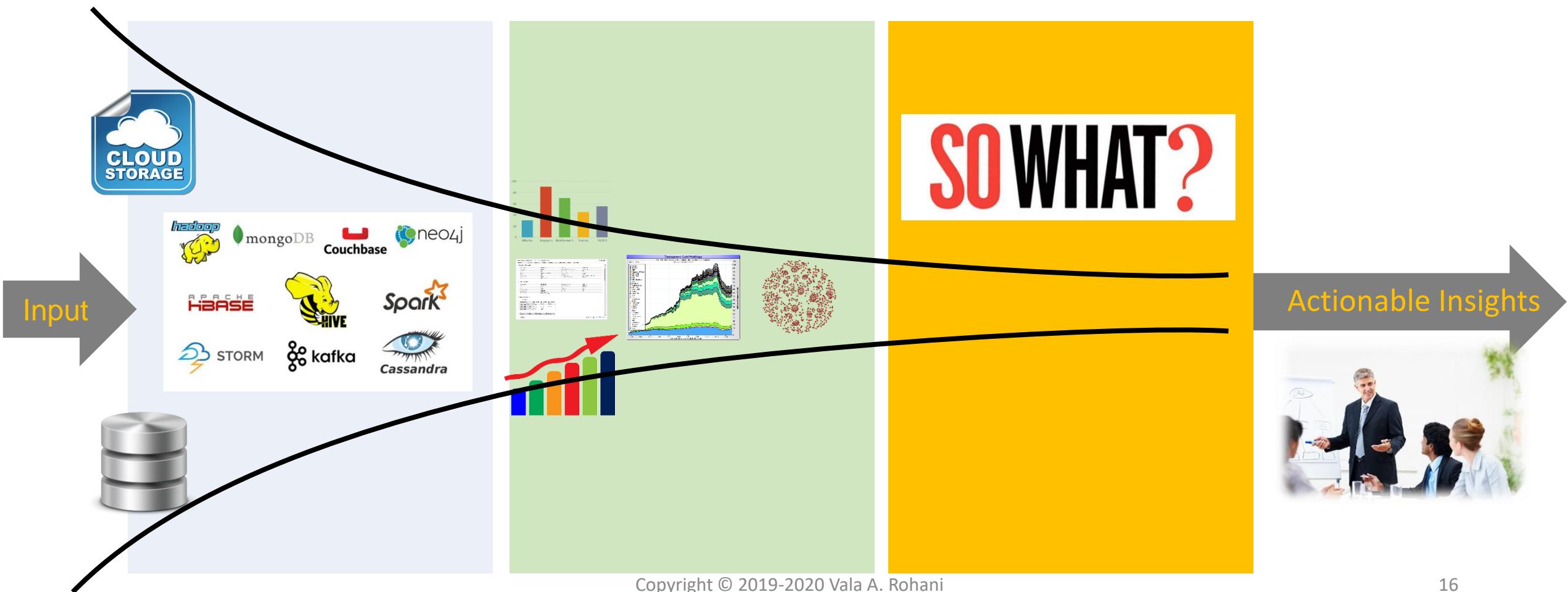


The journey from Data to Actionable Insights

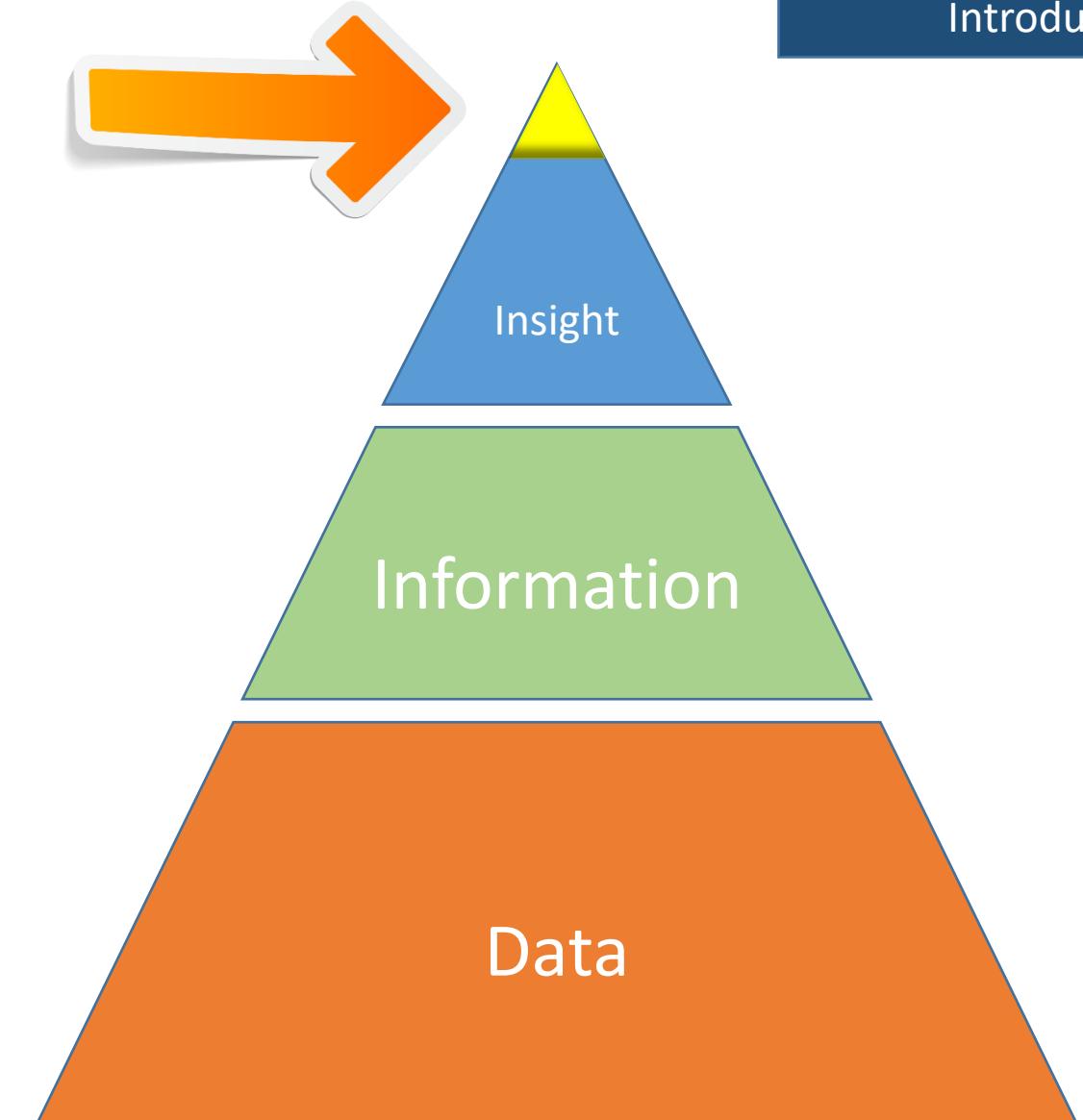
Data

Information

Insights



Actionable Insights



Do not get trapped in your Data and Reports

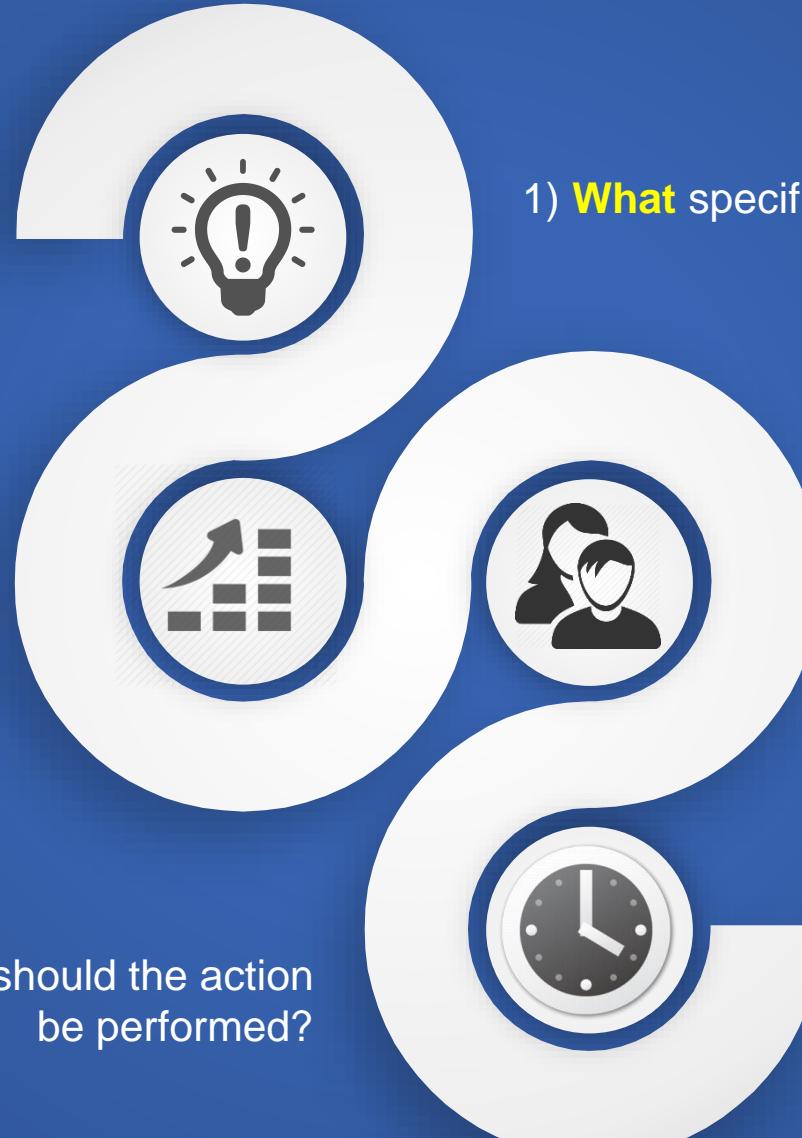
Be an Analytics Ninja because ...



**Actions drives business forward
not data or reports!**



4 Angles of Actionable Insights



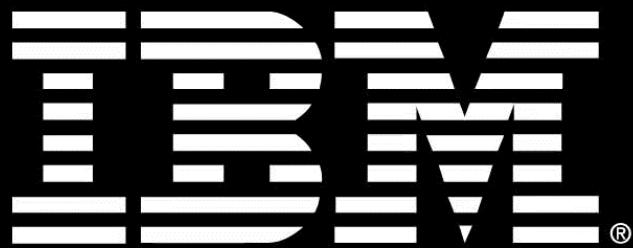
2) **How** will the specific action improve the business

1) **What** specific action to perform?

3) **Who** is responsible for the action?

4) By **when** should the action be performed?

How to generate Actionable Insights?



Data Democratization

Course: Data Analytics
Introduction



Data Democratization

Problem Statement:

- In the past, data has been primarily entrusted to only two privileged groups : **Executives and Data Specialists**
- But After decentralizing the decision-making and increasing responsiveness, Organizations are seeking to empower more workers by granting them the access to related Data
- However, managers have some fear of how people will use and interpret the data.
- As a result, more business questions end up funneling through the analytics team for answers.
- **Unfortunately, analytics teams often can't scale sufficiently to handle the increasing volume of data-related questions**

So what's an aspiring data-driven organization to do in this situation?

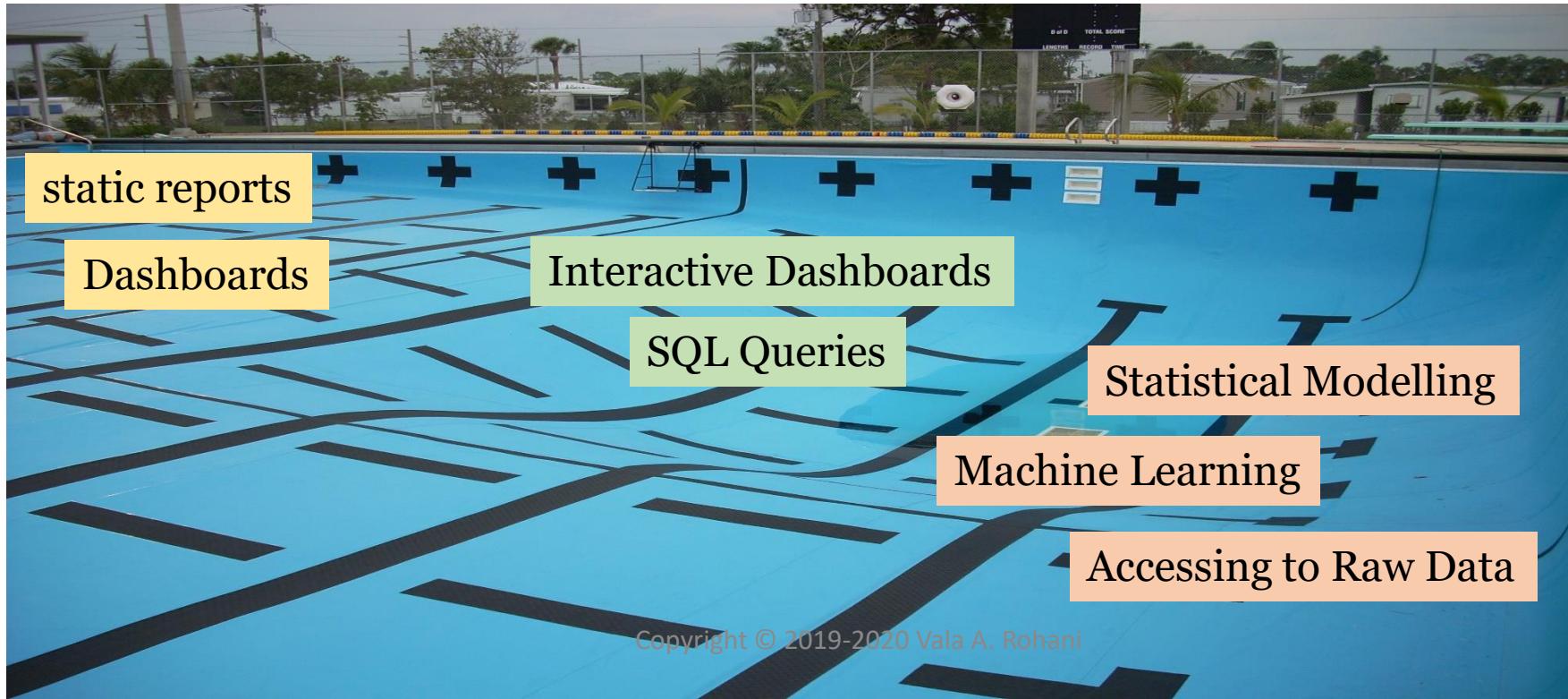
You want your employees to be able to access the data they need, but you also don't want them to harm themselves with it.

Your corporate data is like a giant swimming pool

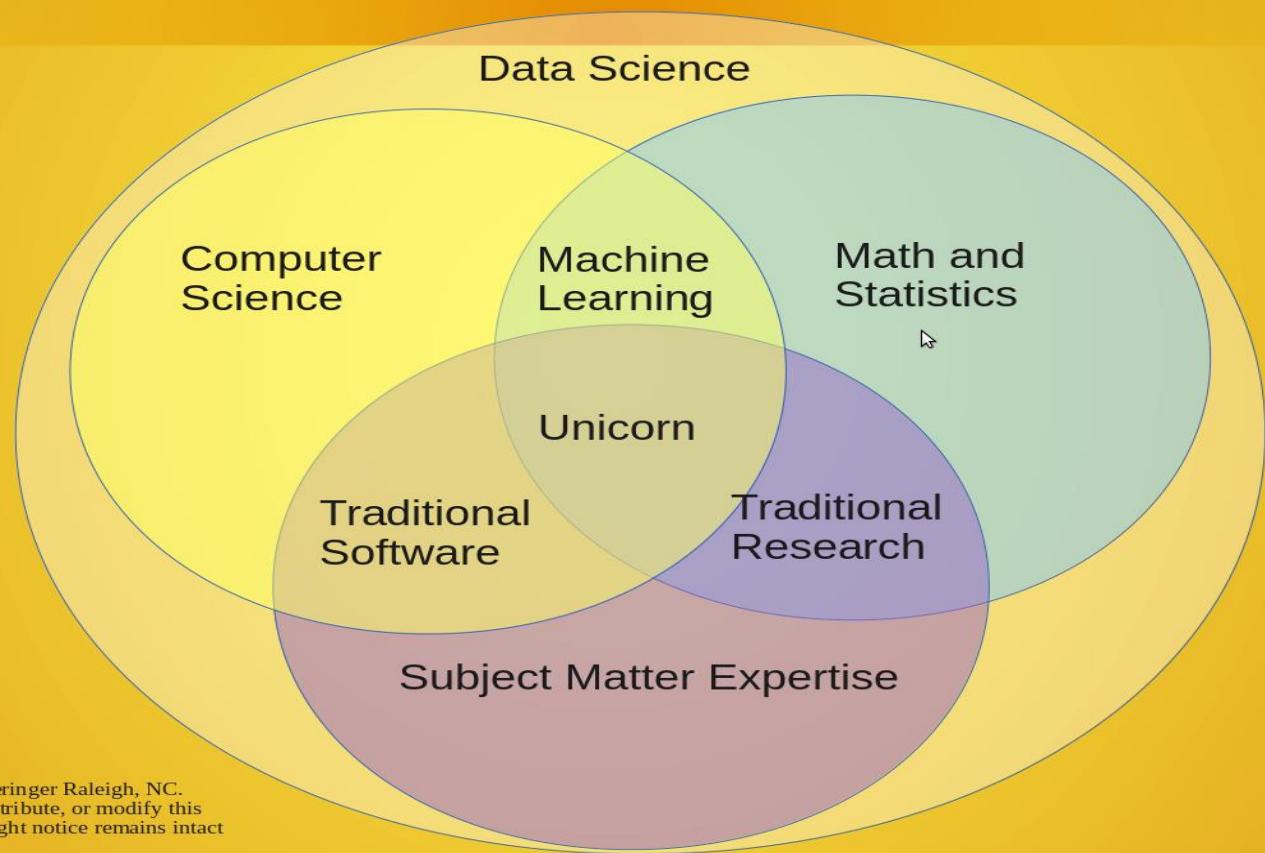
Data Democratization

As the Solution:

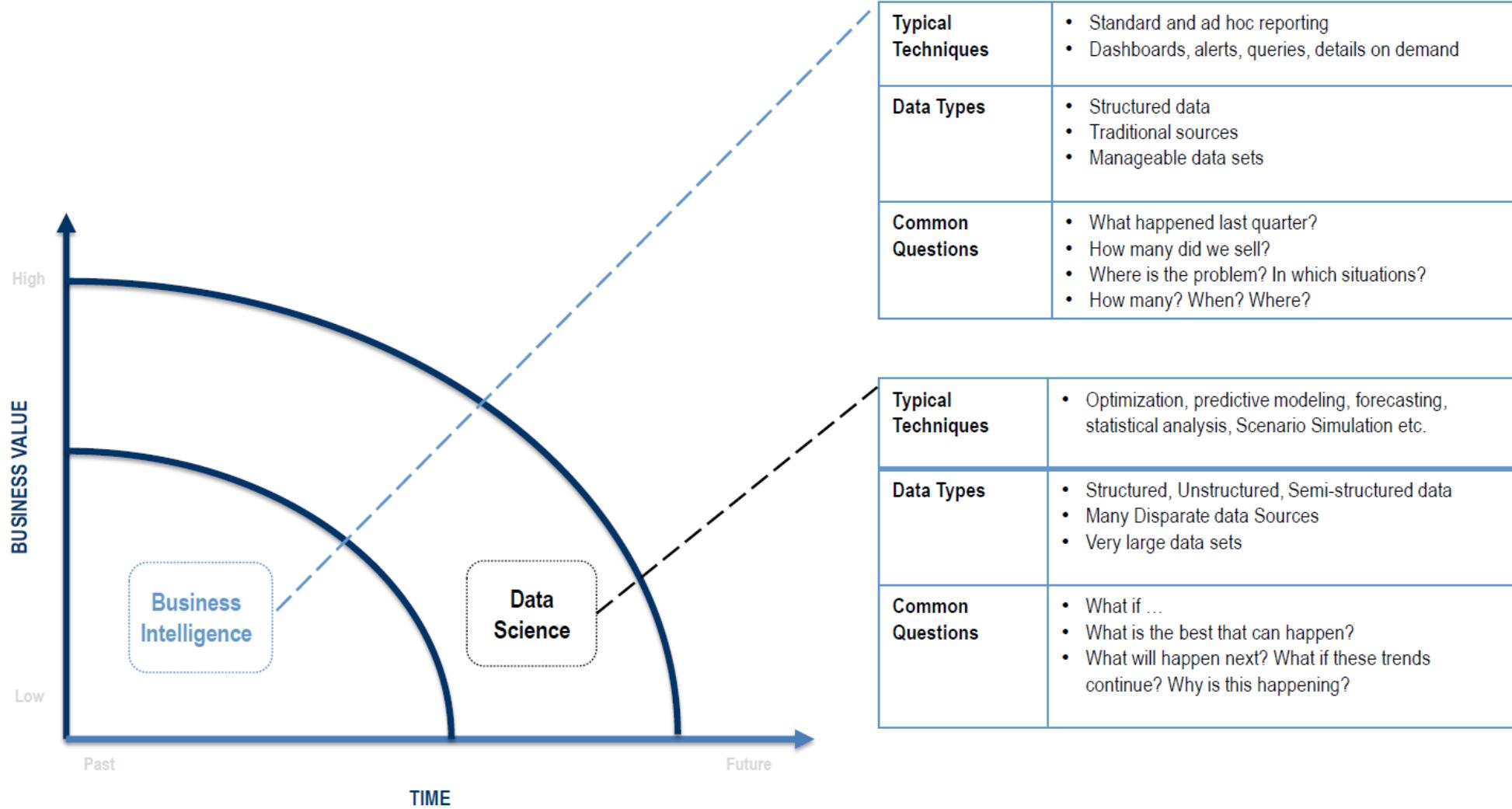
A multi-tiered data analytics approach is recommended, so various users can dive into the right depth of data based on their needs and analytical skills.



Data Science Venn Diagram v2.0



Business Intelligence vs Data Science



AI vs ML vs DL

Artificial Intelligence

A technique for incorporating human intelligence to machine

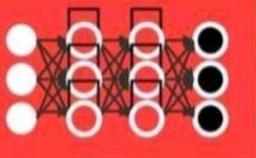
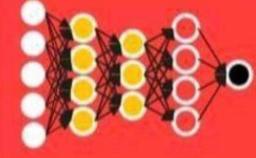
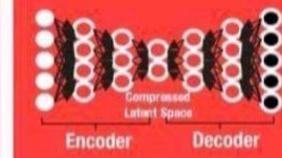
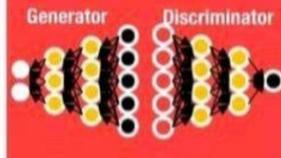
Machine Learning

ML is a subset of AI that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. ML is about **learn from past to predict the future.**

Deep Learning

DL is a subset of ML where **artificial neural networks**, algorithms inspired by the human brain, learn from large amounts of data.

DEEP LEARNING USE CASES

			
RNN	CNN	AUTOENCODER	GAN
 Stock Price Prediction	 Face Recognition	 Dimensionality Reduction	 Generating Image Data
 Chatbots	 Medical Image Detection	 Image Compression	 Generating Art
 Voice Assistant	 Object Detection	 Feature Extraction	 Image-to-Emoji Conversion
 Weather Forecasting	 Character Recognition	 Anomaly Detection	 Face Attribute Manipulation
 Music Generation	 Document Classification	 Fraud Detection	 3D Object Generation



DDO?

Data Driven Organization

What is **NOT** DDO?



Having Lots of reports does not make you data driven!

What is **NOT** DDO?



Having Lots of dashboards does not make you data driven!

What is **NOT** DDO?



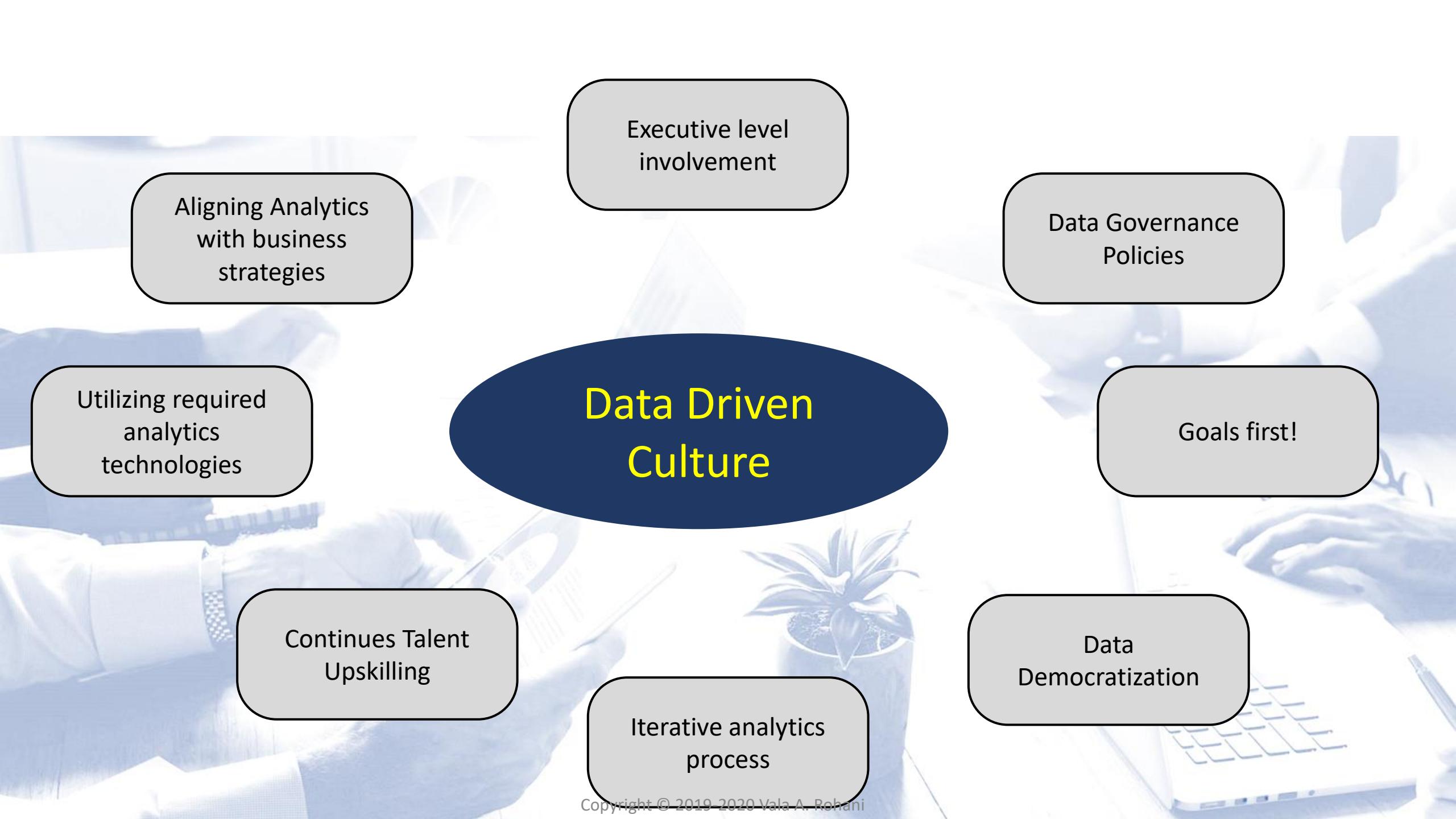
Having **BIG** Data clusters does not make you data driven!

Copyright © 2019-2020 Vala A. Rohani

What **is** DDO?

To become DDO, you need to create a **Data-Driven Culture** in your organization





Data Driven Culture

Aligning Analytics
with business
strategies

Executive level
involvement

Data Governance
Policies

Utilizing required
analytics
technologies

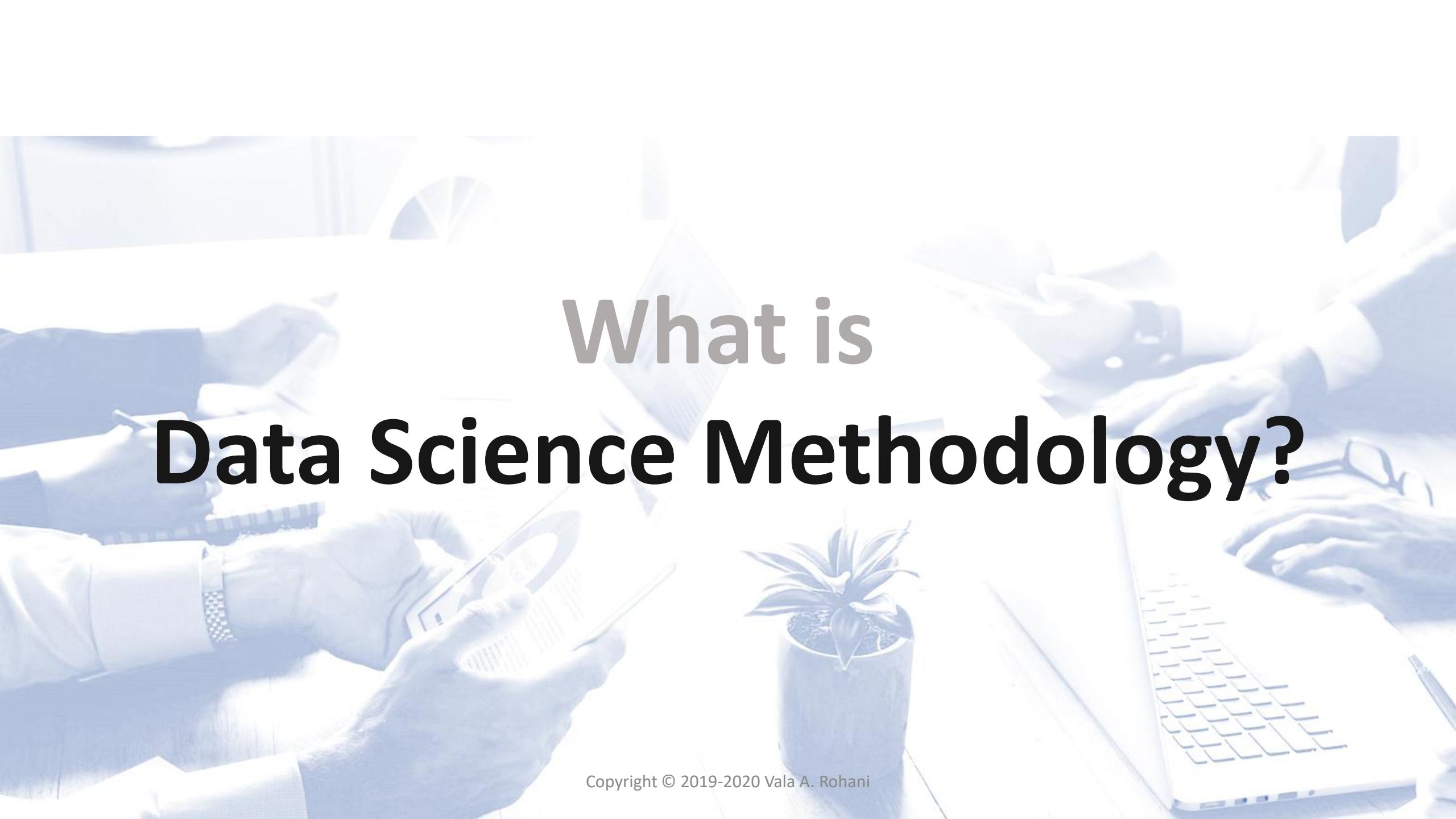
Goals first!

Continues Talent
Upskilling

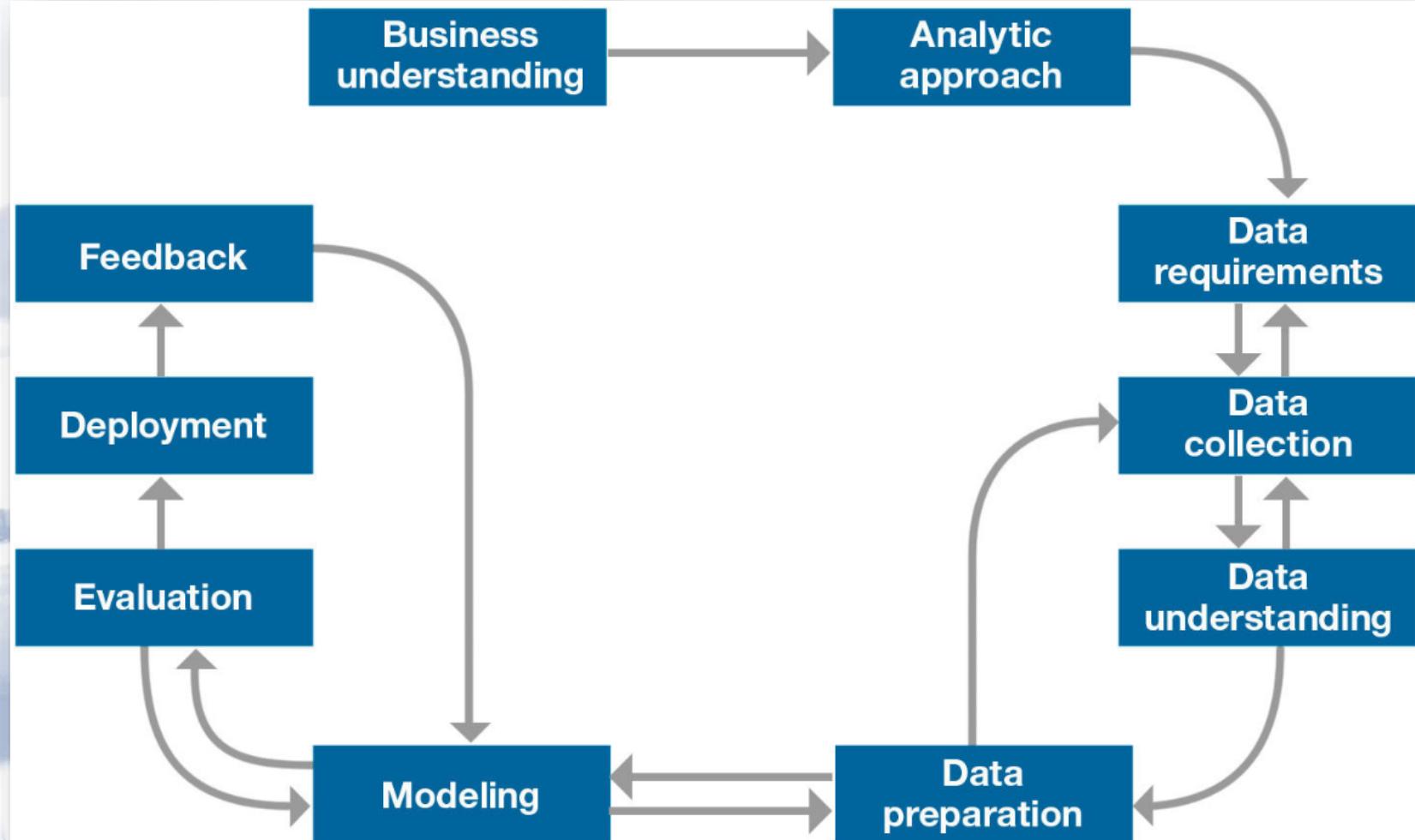
Data
Democratization

Iterative analytics
process

What is **Data Science Methodology?**



Data Science Methodology



<https://www.ibmbigdatahub.com/blog/why-we-need-methodology-data-science>



What skills you need to become Data Scientist?

DATA SCIENCE SKILL SET

BASIS

CORE

ENABLEMENT

TECHNICAL SKILLS

- CODING SKILLS
- HANDLING DATA
- COMPUTATIONAL TOOLS
- BASIC SOFTWARE DEVELOPMENT
- BIG DATA
- HIGH PERFORMANCE COMPUTING
- PARALLEL COMPUTING

TECH-SAVVY

ANALYTICAL SKILLS

- ADVANCED STATISTICS & INFERENCE
- MODELLING & SIMULATION
- MACHINE LEARNING
- COMPUTER SCIENCE
- ADVANCED MATH
- DATA VISUALIZATION
- EXPERIMENT DESIGN
- RESEARCH EXPERTISE

SCIENTIFIC PROBLEM SOLVING

BUSINESS SKILLS

- EVALUATION AND DEVELOPMENT OF BUSINESS CASES
- PROJECT MANAGEMENT
- BUSINESS PROCESSES
- CHANGE MANAGEMENT
- COMMUNICATION SKILLS
- LEADERSHIP SKILLS

APPETITE FOR BUSINESS PROBLEMS



Further details?!



Data Science: A mix of hard and soft skills



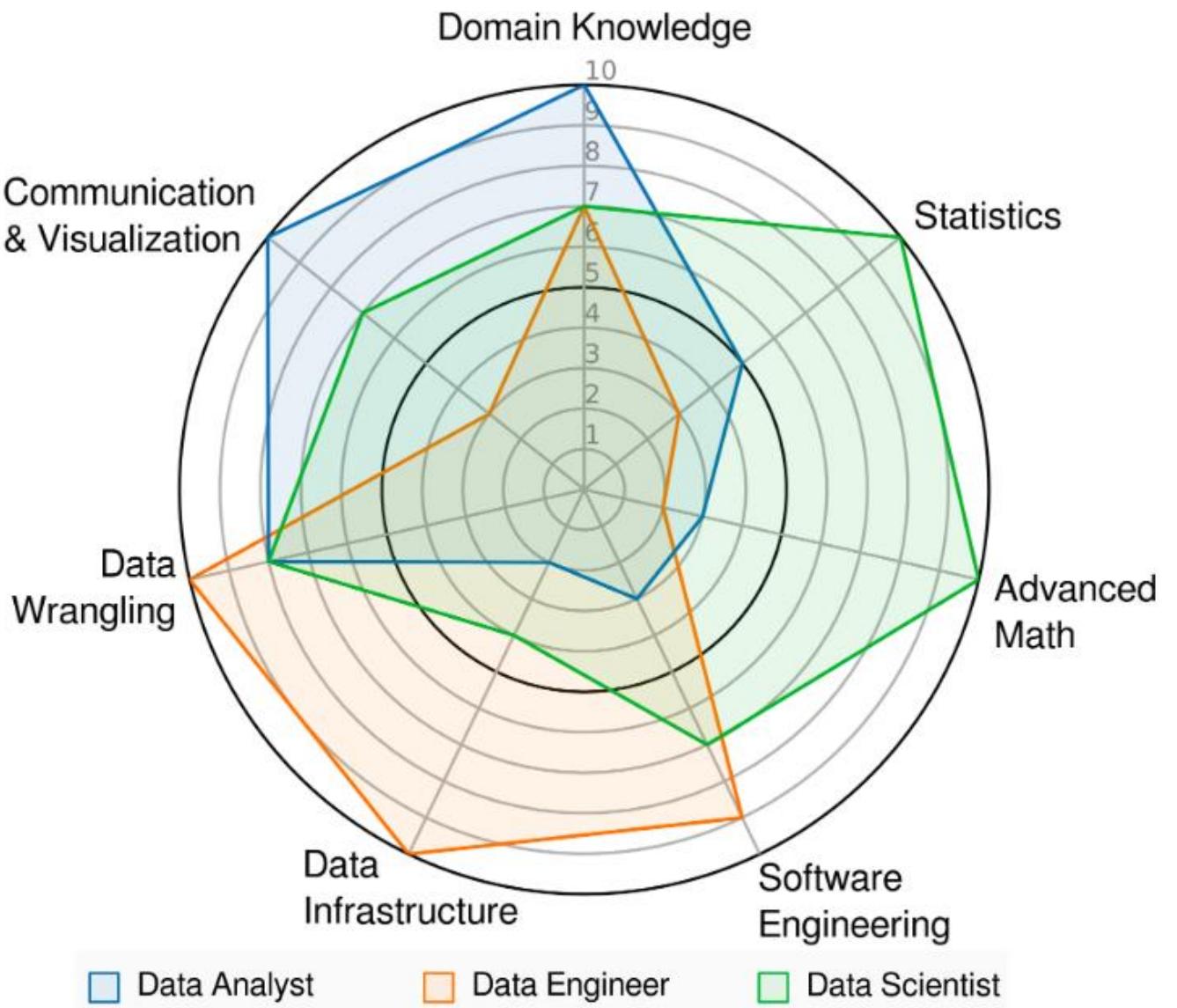
Data science is one area of the digital sector that is desperately in short of talent. According to [EDM Monitoring Tool](#):

Europe will face a shortage of **925,000** data professionals by 2025.



The Data Science Team: Who Does What





How to accelerate the adoption of **data science** within an organization?



How to accelerating the adoption of *data science* within an organization?

1. Access to Data, in the rawest form
2. Having strong data governance practices in place
3. Provide an analytic sandbox
4. Leverage data visualization tools
5. Keep your team structures flexible
6. Integrate data scientists with existing business units
7. Be sure your data science sponsors are 100% committed to supporting the team

<http://www.informationweek.com/big-data/big-data-analytics/how-to-build-a-successful-data-science-team-/d/d-id/1113234>

Prepared and Presented by Dr Vala Ali Rohani



How to accelerating the adoption of *data science* within an organization?

8. Develop internal team skillsets
9. Iterate quickly
10. Be comfortable with data not being perfect
11. Plan now for scaling up
12. Think about producing Actionable Insights



How to manage a data science project



Prepared and Presented by Dr Vala Ali Rohani

What is your main responsibility as
Data Scientist?

*To use Data and Analytics to solve
stakeholder business problems.*





To this end, we need to know ...

What the real business problem is?

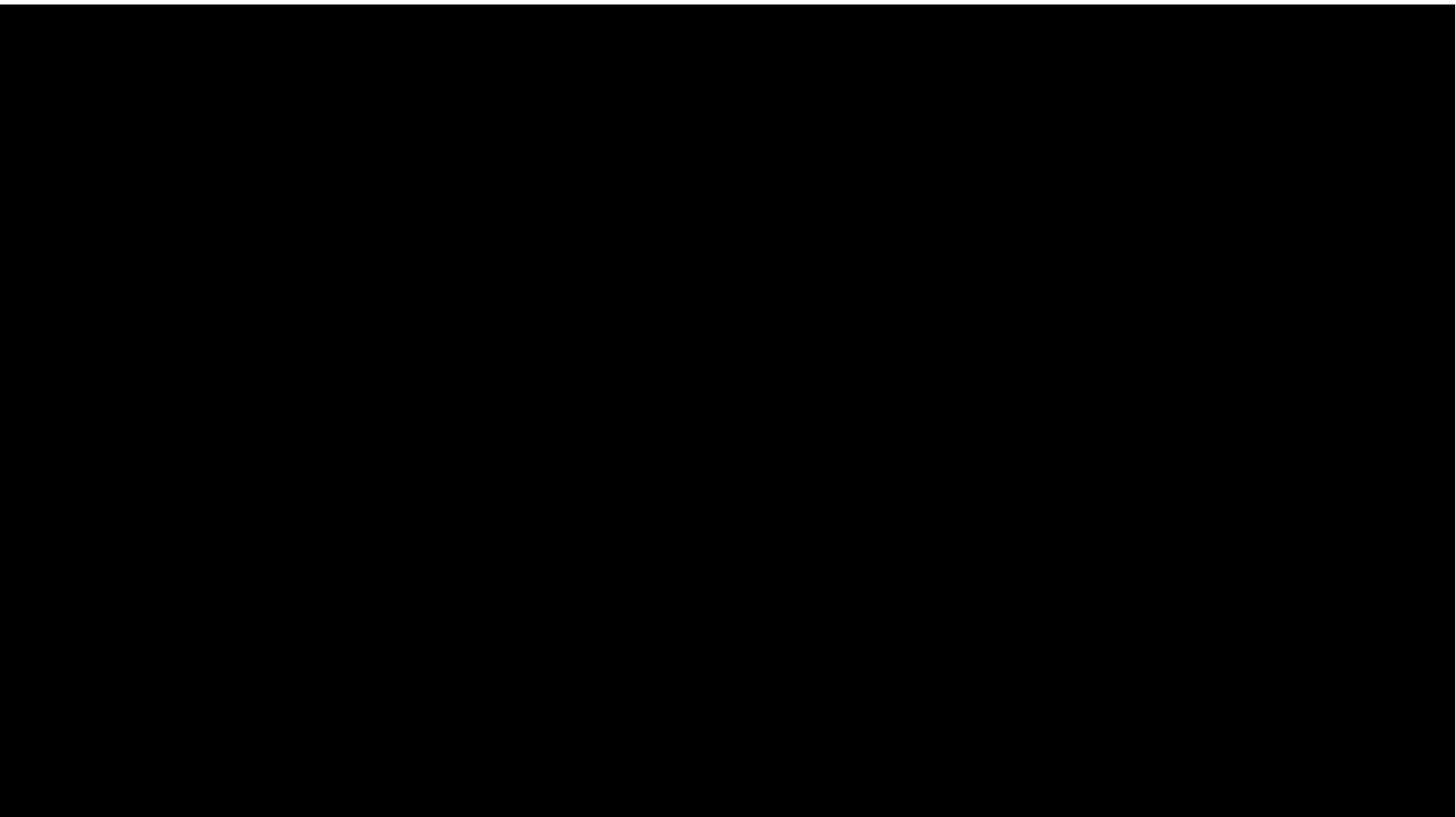
And, ask questions ...

Meet the primary people: Stakeholders



- What the problem is in this business that you are hoping to solve using analytics?
- How this problem is affecting your business?
- What is your ideal outcome for this project?

Let's watch a story about one of real challenges in Analytics Projects!



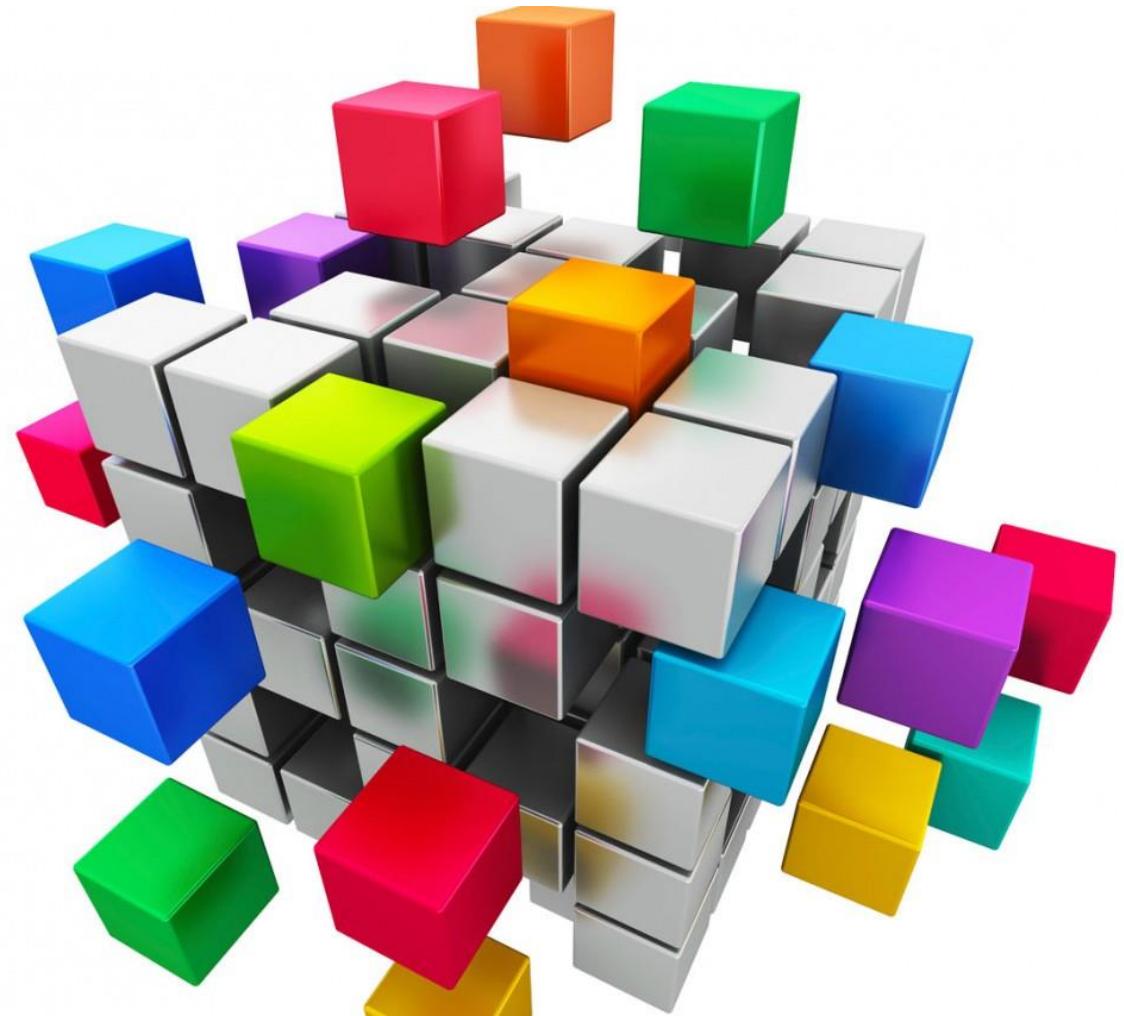
To

Identify effective objectives

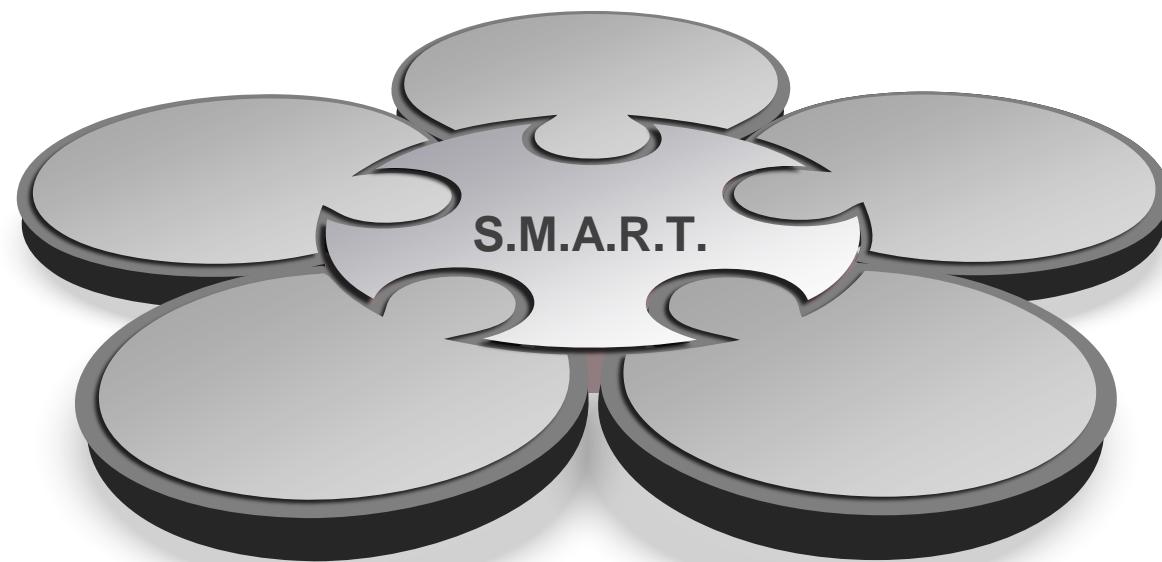
But How???

S.M.A.R.T. Framework

George Duran, 1981



The establishment of projects objectives is suggested be created using the **S.M.A.R.T.** philosophy.



Let's start with an example:

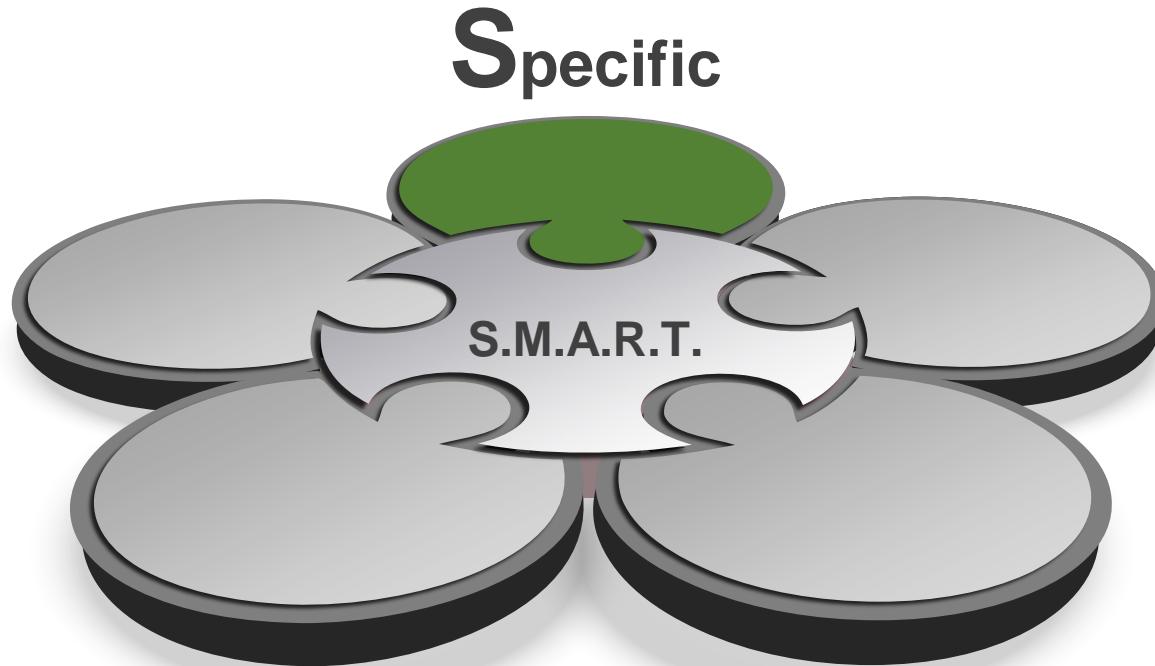
Imagine your client is an online retailer of the men's clothing store.

Project manager describe the problem as follows:

People aren't returning back to our website after they login for the first time. I need you to tell us how convert first time visitors to returning visitors.

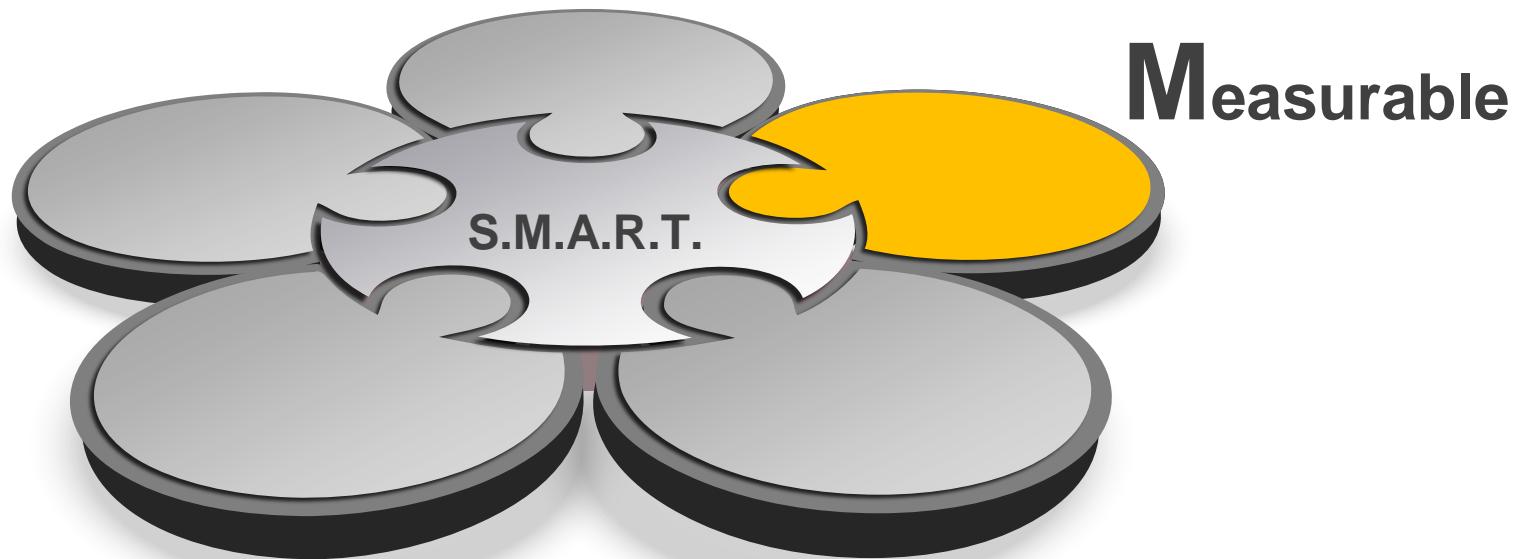


Phrasing the Project Goal using **SMART** method:



Increasing the number of returning visitors

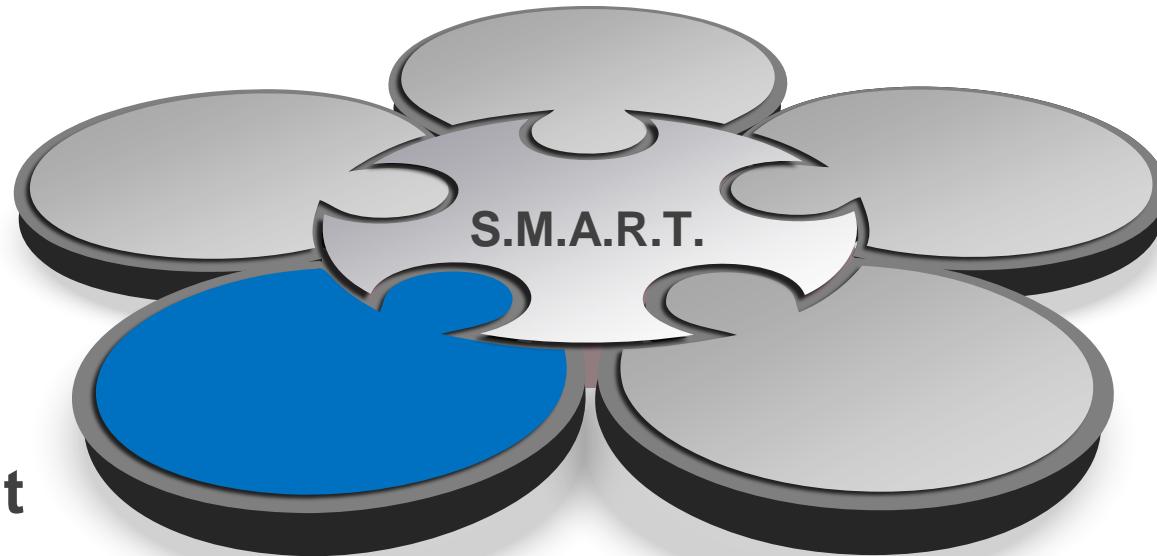
Phrasing the Project Goal using **SMART** method:



Increasing the number of returning visitors **on a month-by-month basis by at least 15% compare to the last month during last year.**

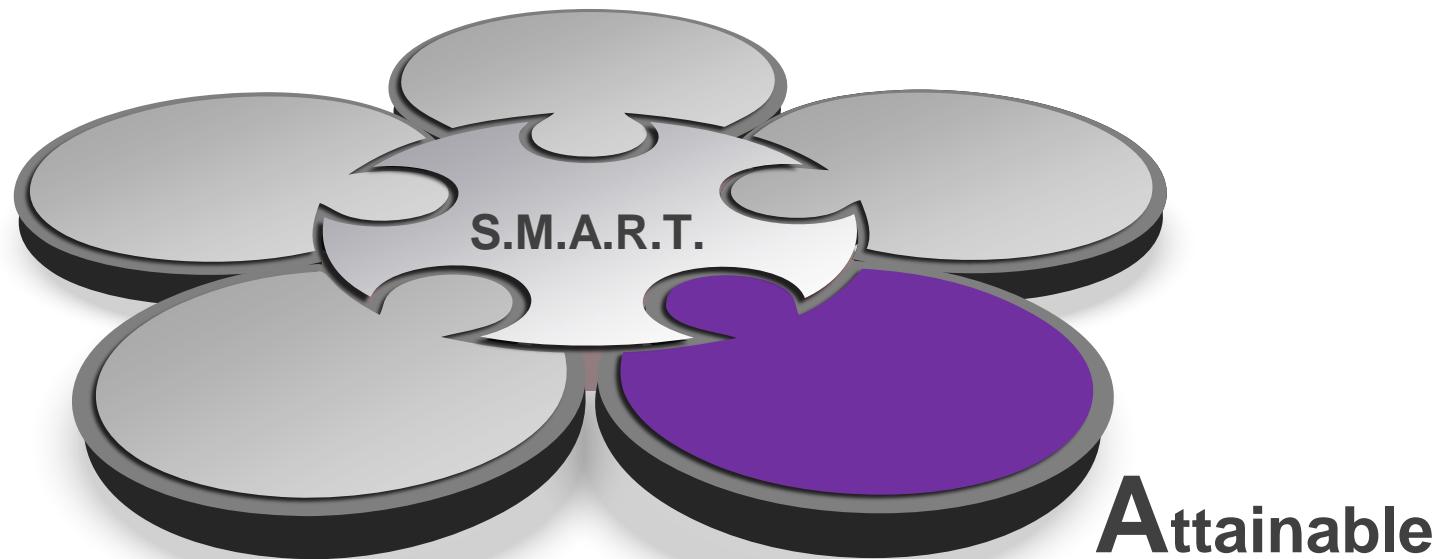
Phrasing the Project Goal using **SMART** method:

Relevant



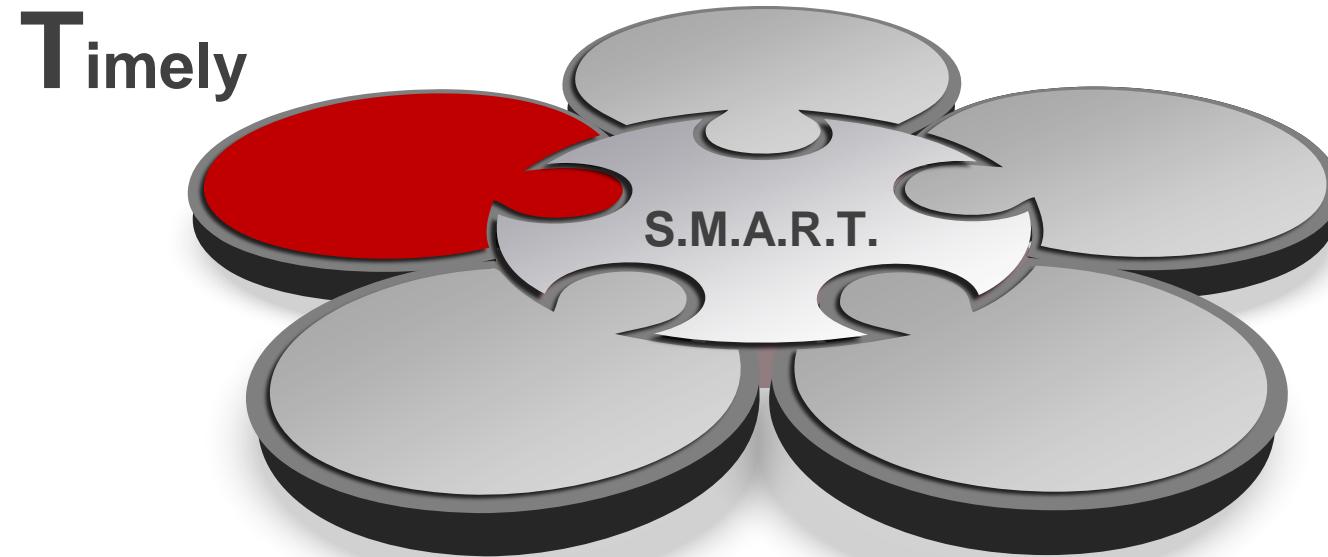
The goal of the project is **to determine the website changes that will most efficiently increase revenues** by 15% on a month-by-month basis compared to the same month last year.

Phrasing the Project Goal using **SMART** method:



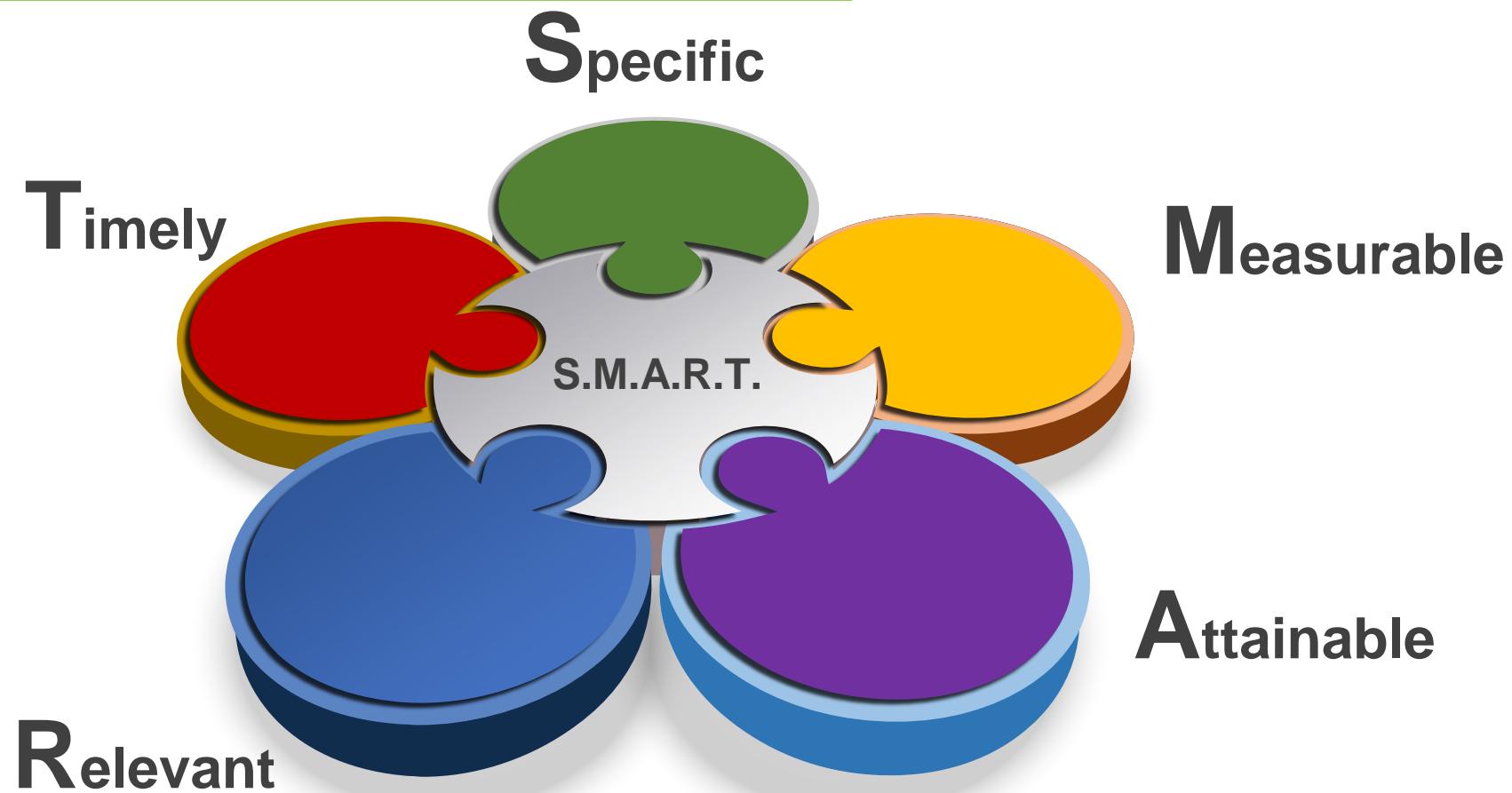
The goal of the project is **to analyze archived click-stream data** to determine the website changes that will most efficiently increase revenues by 15% on a month-by-month basis compared to the same month last year.

Phrasing the Project Goal using **SMART** method:



The goal of the project is to, **by the end of two months**, analyze archived click-stream data to determine the website changes that will most efficiently increase revenues by 15% on a month-by-month basis compared to the same month last year.

Phrasing the Project Goal using **SMART** method:



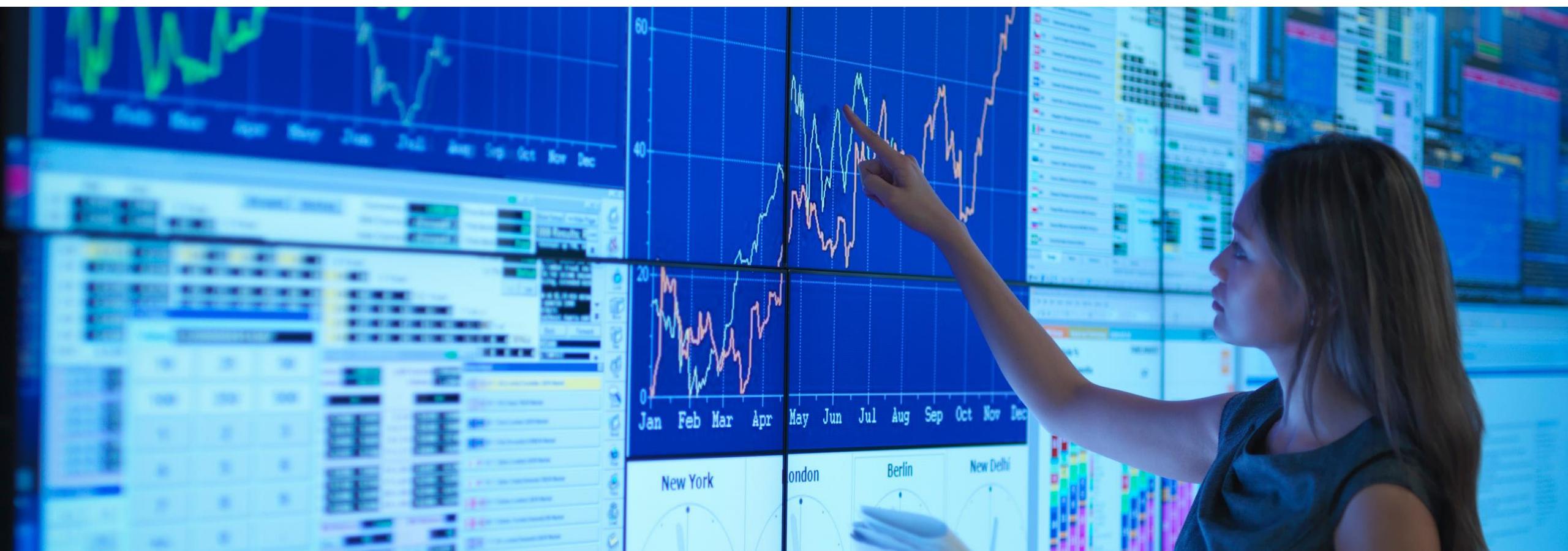
Phrasing the Project Goal using **SMART** method:

Here's a more complete version, that might have two parts:

First in 3 months, install a system that will collect and store click-stream data in a cloud-base relational database.

Then, the second part could be. By the end of 2 months, after the system has been installed, this data will be analyzed to determine the website changes that will most efficiently increase revenues by 15% on a month-by-month basis compared to the same month last year.

How to Design a Data Analytics Project?



Structured Pyramid Analysis Plan

S
P
A
P

Question: What's a variable?

Answer: A variable is an object, event, idea, feeling, time period, or any other type of category you are trying to measure. There are two types of variables-independent and dependent.

S
P
A
P

Question: What's an independent variable?

Answer: It is a variable that stands alone and isn't changed by the other variables you are trying to measure.

For example, someone's age might be an independent variable.

In fact, when you are looking for some kind of relationship between variables you are trying to see if the independent variable causes some kind of change in the other variables, or dependent variables.

S

P

A

P

Question: What's a dependent variable?

Answer: It is something that depends on other factors.

For example, a test score could be a dependent variable because it could change depending on several factors such as how much you studied, how much sleep you got the night before you took the test, or even how hungry you were when you took it.

Usually when you are looking for a relationship between two things you are trying to find out what makes the dependent variable change the way it does.

S
P
A
P

Some Examples:

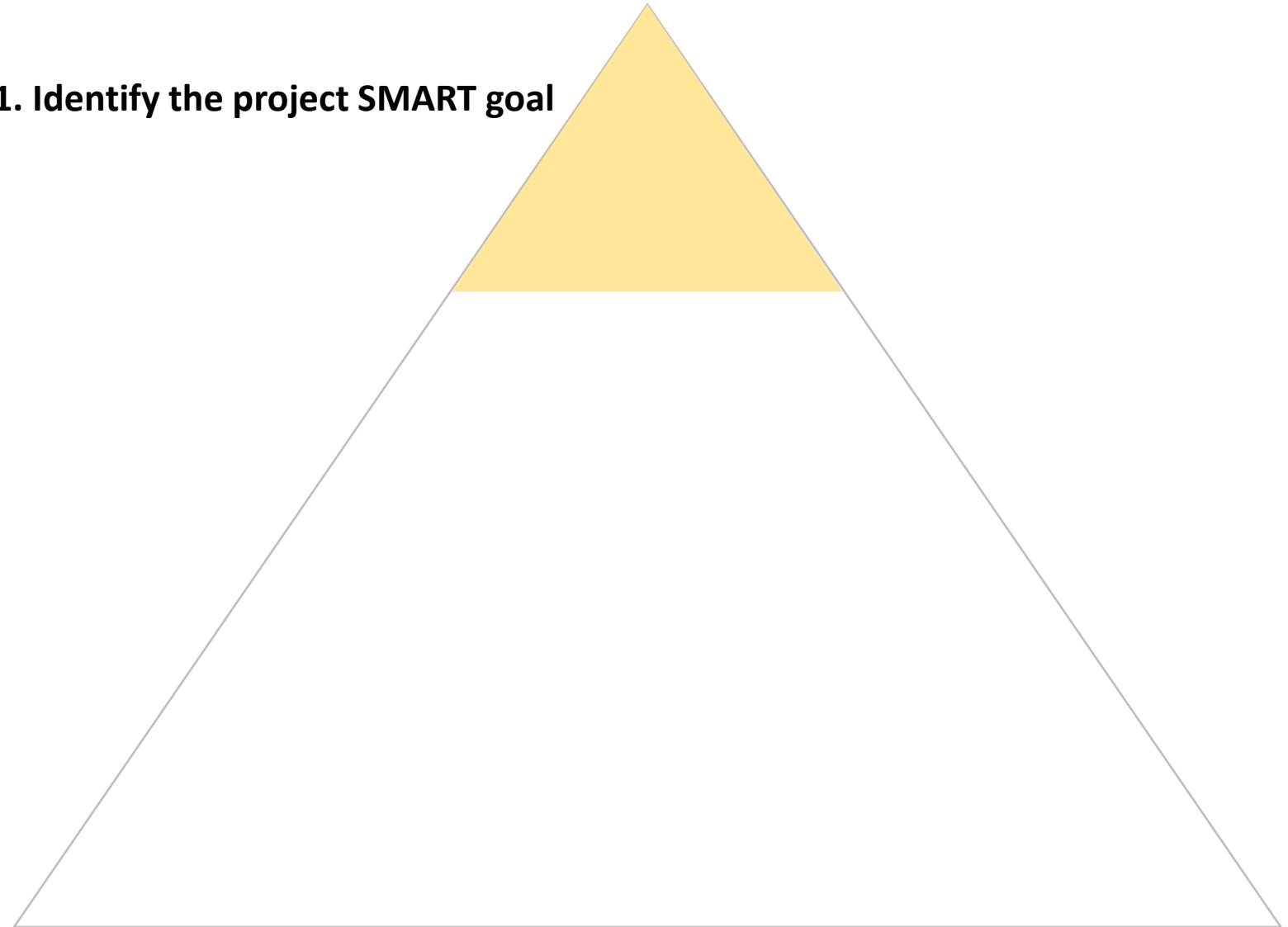
1. You are interested in whether a higher minimum wage impacts employment rates.
2. A scientist studies the impact of a drug or age or smoking on cancer.
3. Educators are interested in whether participating in after-school math tutoring can increase scores on standardized math exams.

Dependent var

Independent Var

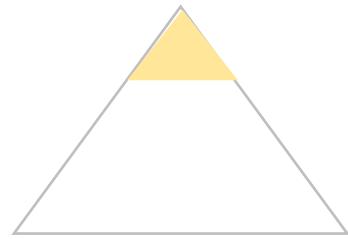
S
P
A
P

1. Identify the project SMART goal



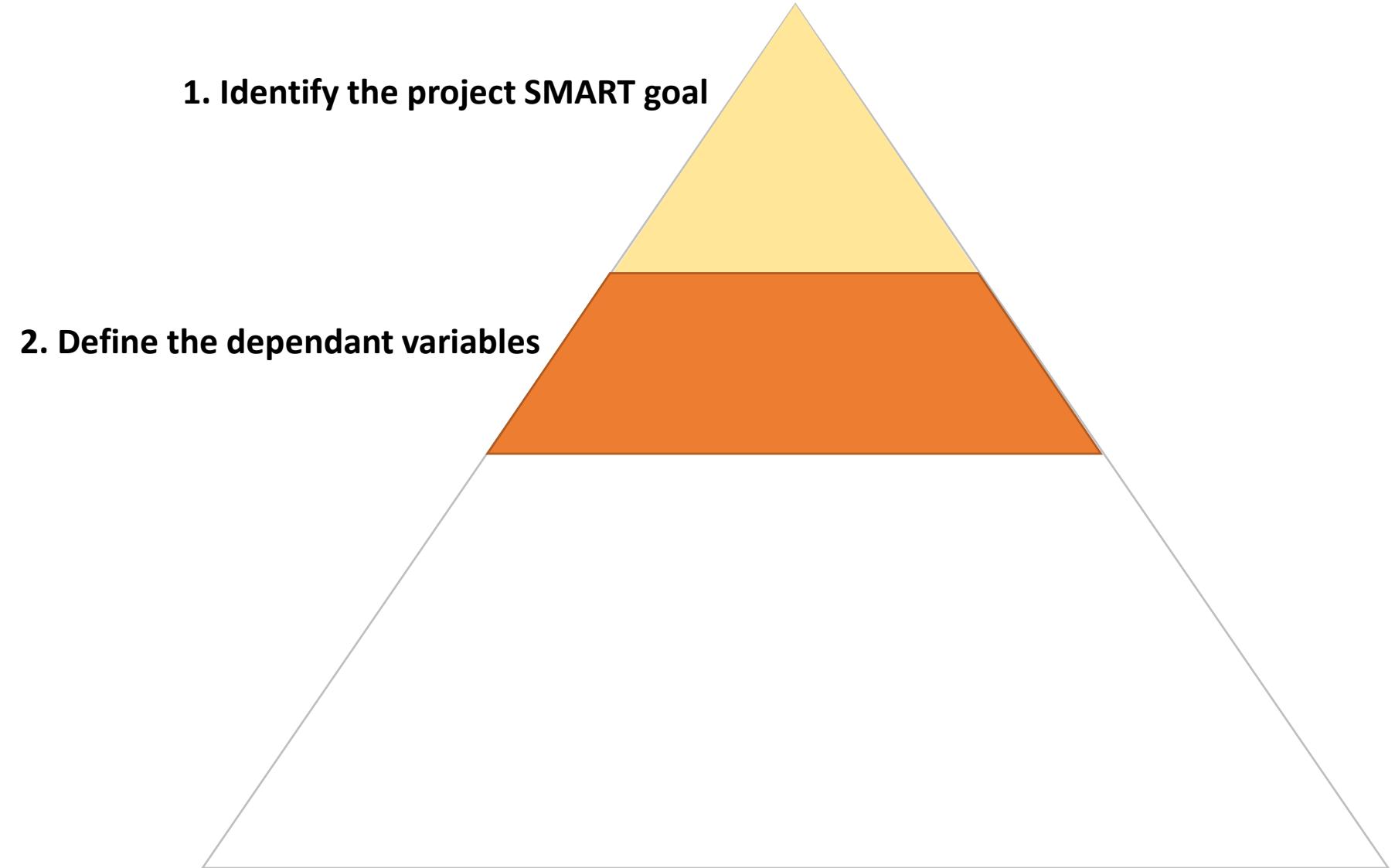
S
P
A
P

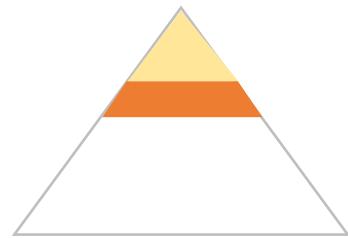
1. Identify the project SMART goal



In 2 months, analyse archived click-stream data to determine the website changes that will most efficiently increase revenues by 15% on a month-to-month basis compared to the same month last year.

S
P
A
P





2. Define the dependant variables

DV1: Total \$ spent per transaction

In clickstream database, “total spent” field aggregated by SUM over each transaction ID

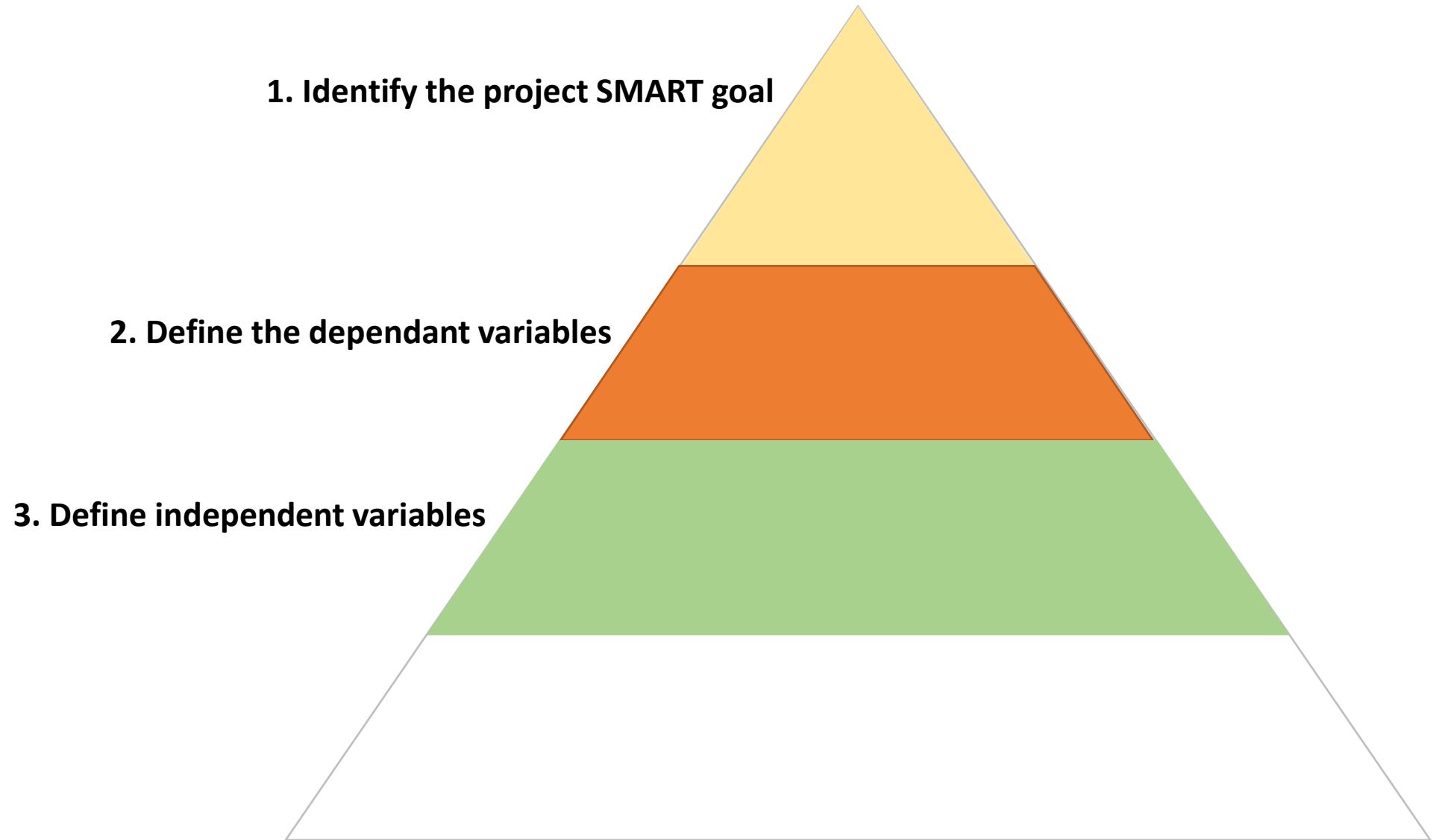
DV2: Total \$ spent per month

In clickstream database, “total spent” field aggregated by SUM over date (month)

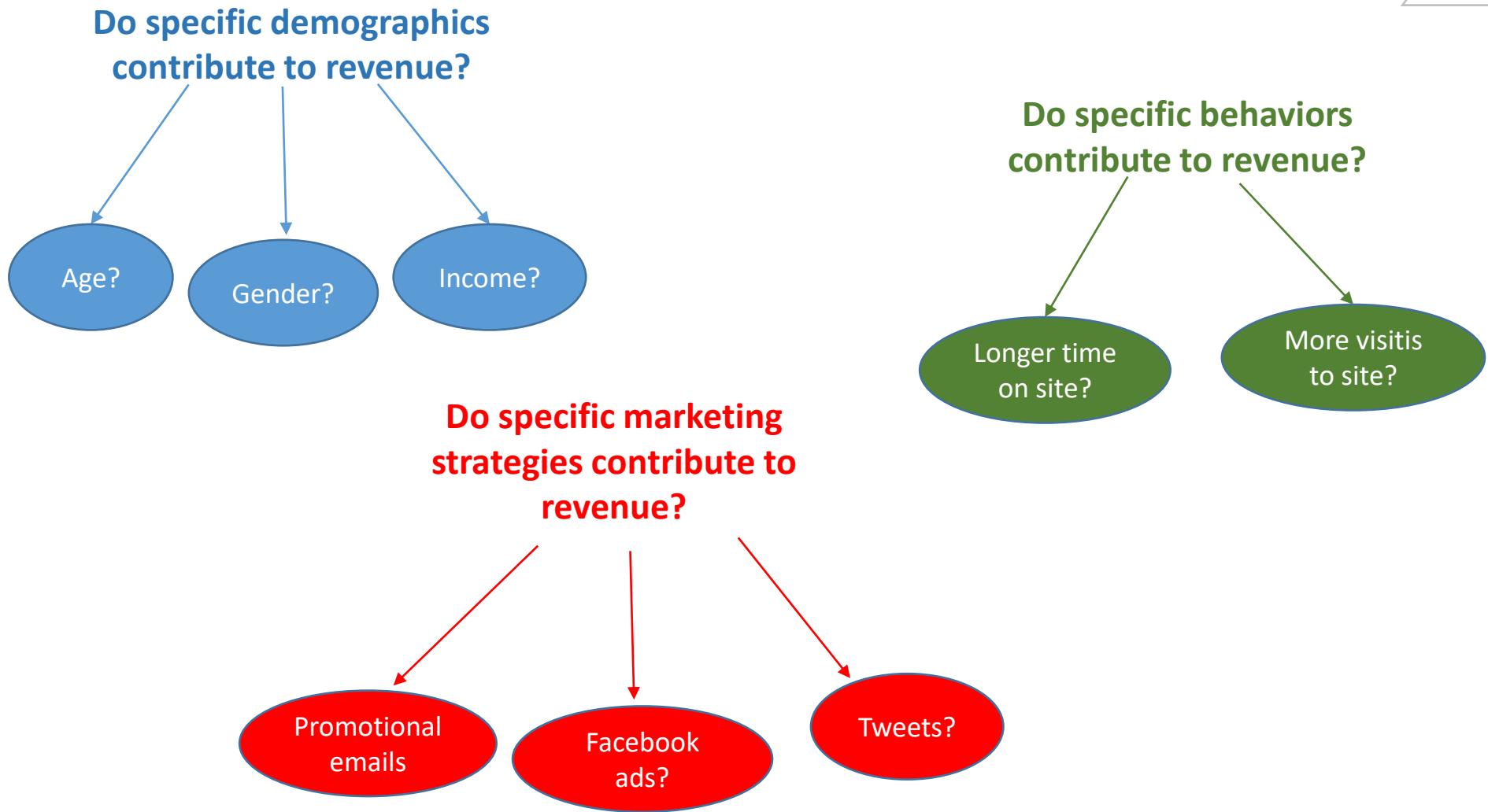
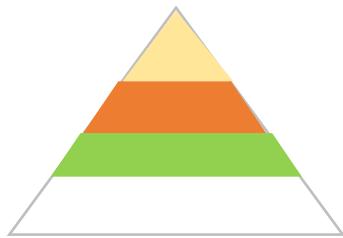
DV3: Total \$ spent per customer

In clickstream database, “total spent” field aggregated by SUM over each customer ID

S
P
A
P

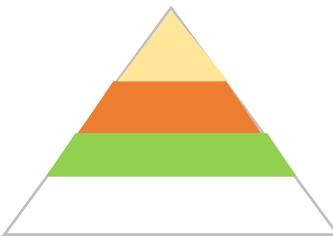
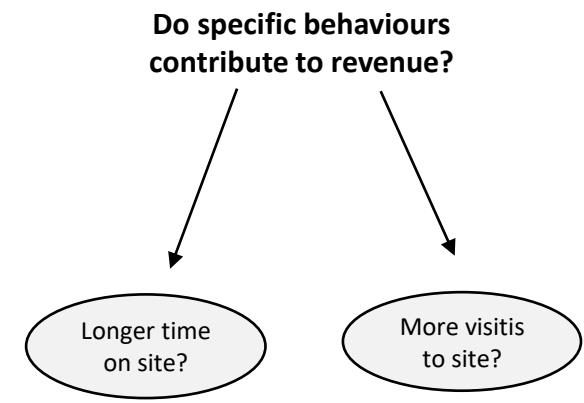
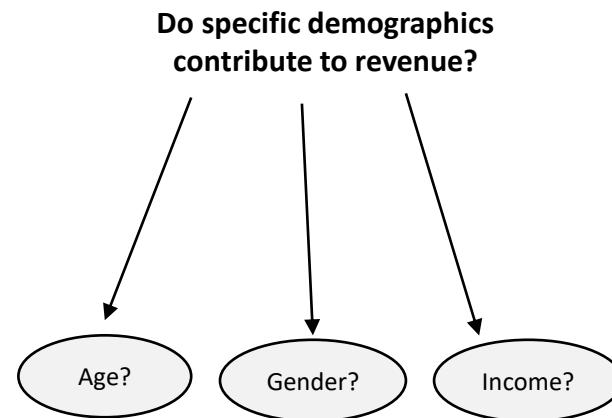


3. Define independent variables

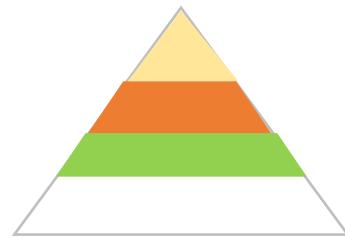


S P A P

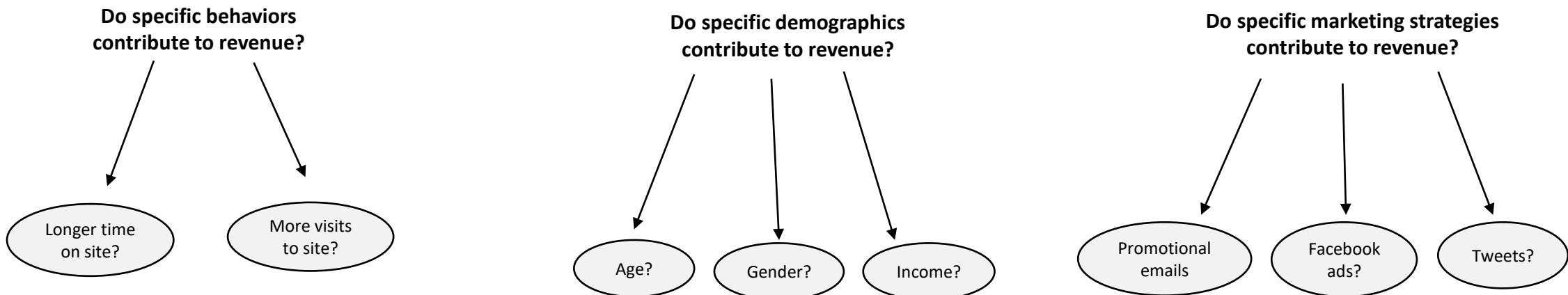
3. Define independent variables



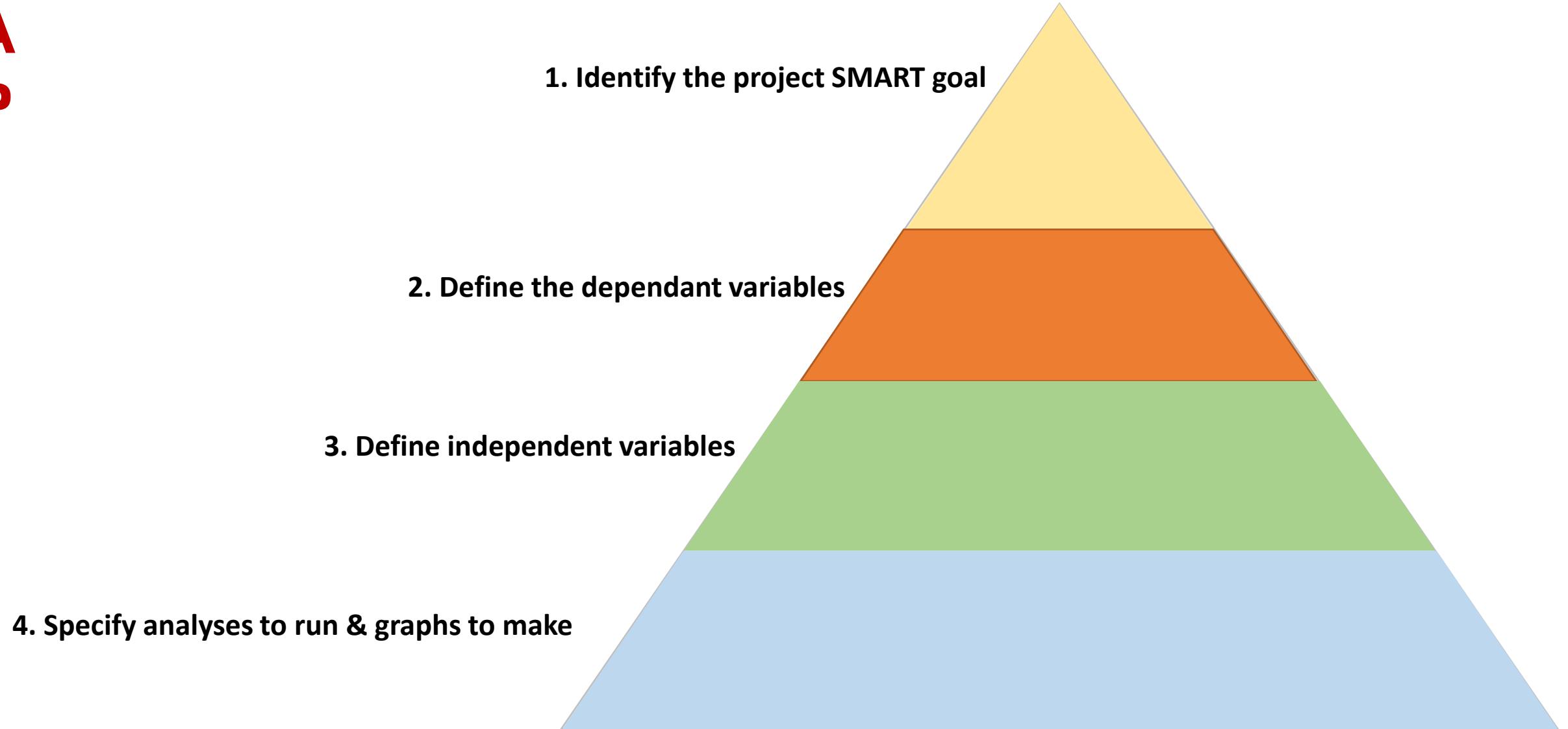
3. Define independent variables



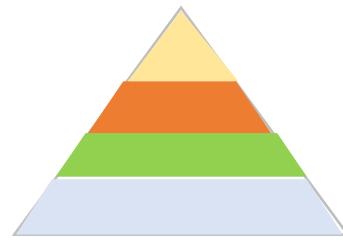
Once you know what data variables you are going to use for each category or subcategory, the next step is to assign each category or subcategory a priority based on **who suggested them, how much impact you suspect they could have, and how feasible you suspect they will be to assess.**



S
P
A
P



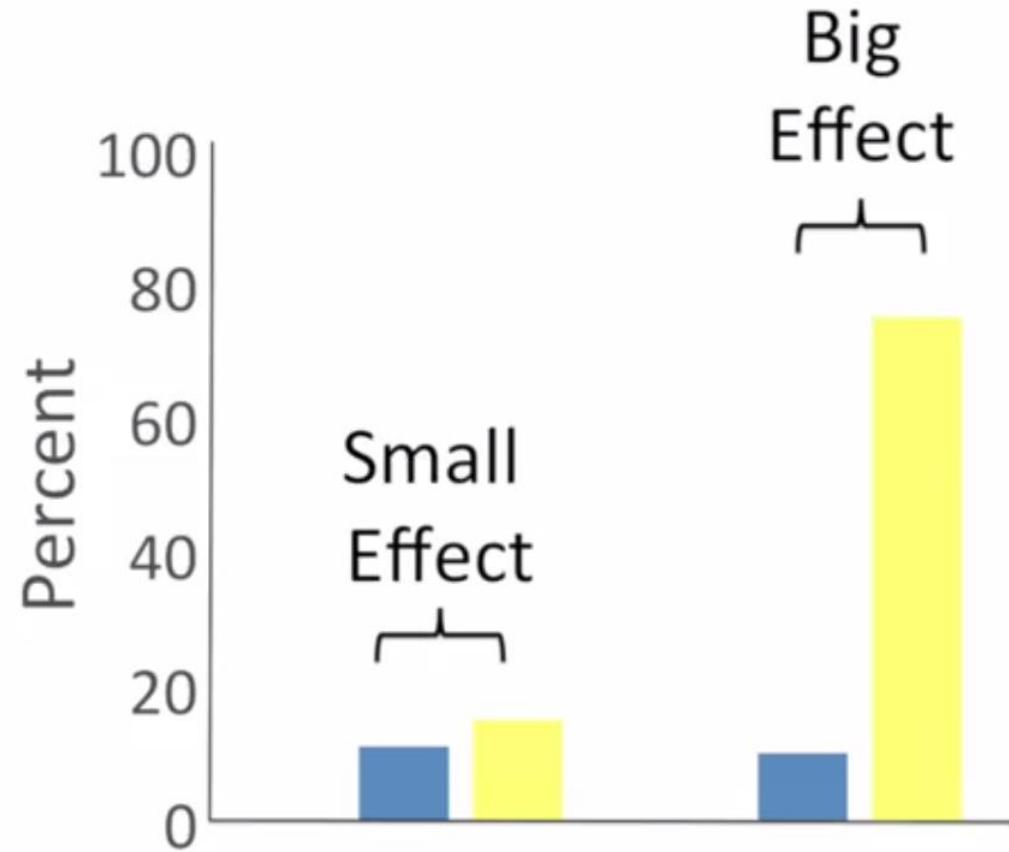
4. Specify analyses to run & graphs to make



The final step of an SPAP is to actually implement the plan and illustrate the most important issues or independent variables that will help you solve your business problem.

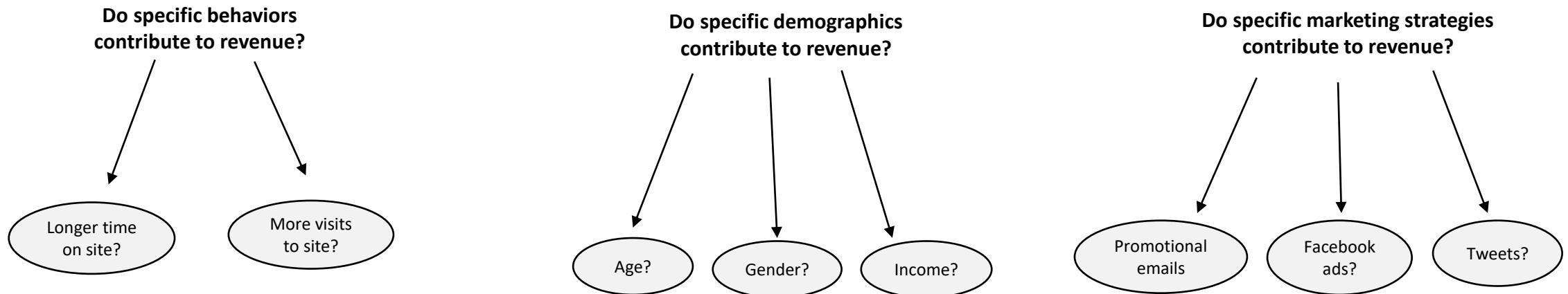
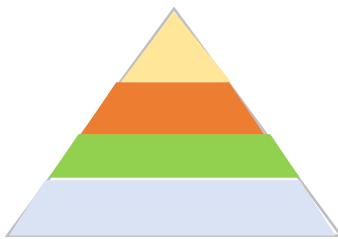
But How?

Using Data Visualizing Software



S P A P

4. Specify analyses to run & graphs to make



Line Chart:

Visit on x-axis
\$ Spend on y-axis

Bar Chart:

Age on x-axis
\$ Spend on y-axis

Bar Chart:

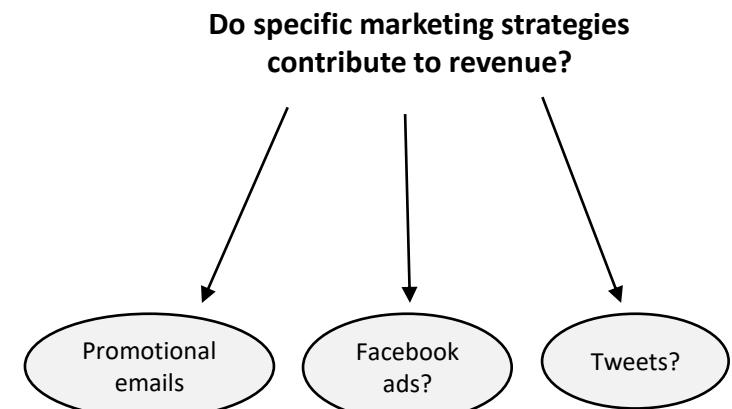
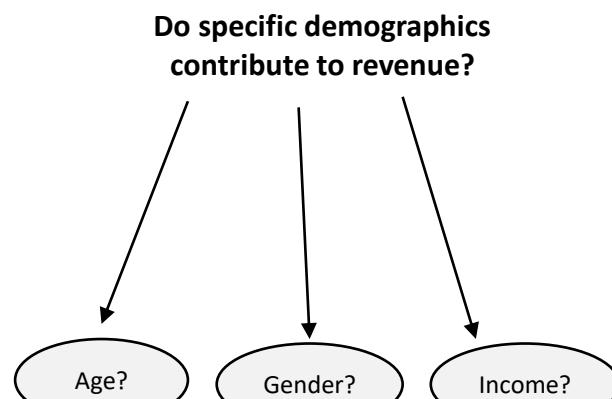
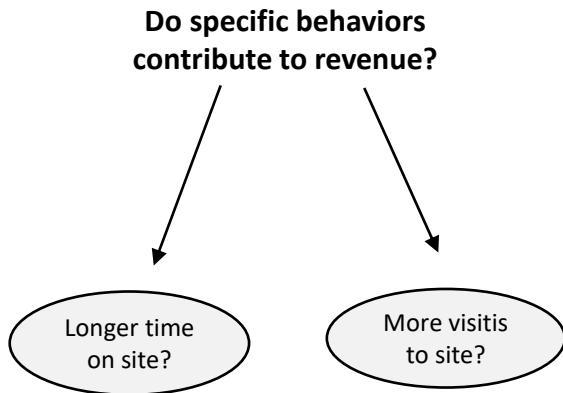
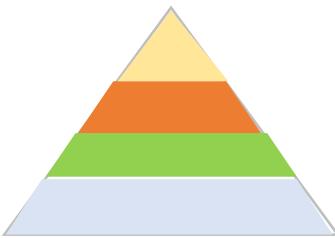
Income on x-axis
\$ Spend on y-axis

Bar Chart:

Gender on x-axis
\$ Spend on y-axis

S P A P

4. Specify analyses to run & graphs to make



Line Chart:

Visit on x-axis
\$ Spend on y-axis



Bar Chart:

Age on x-axis
\$ Spend on y-axis

Bar Chart:

Income on x-axis
\$ Spend on y-axis

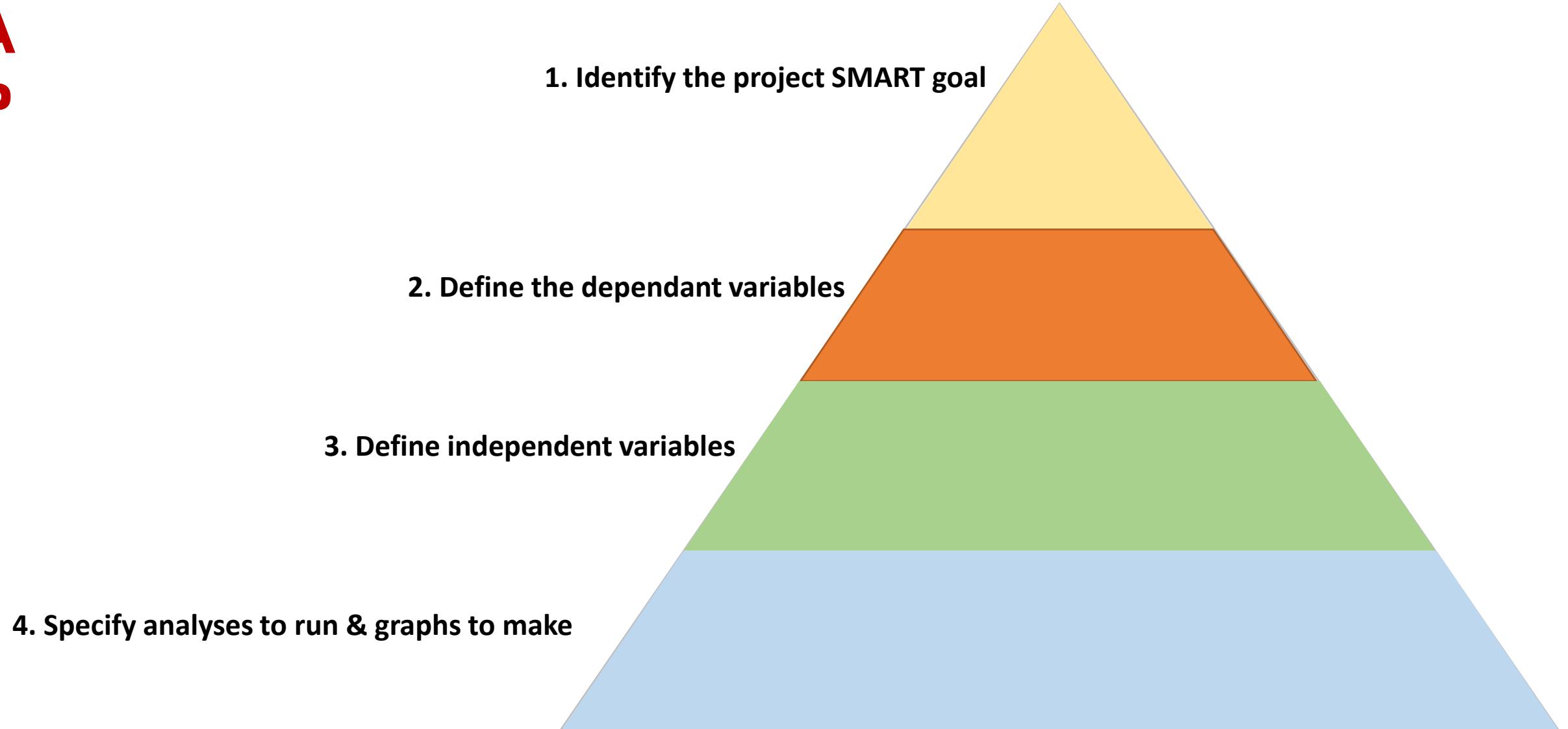


Bar Chart:

Gender on x-axis
\$ Spend on y-axis

Do any eye-catching effect you observe?

S
P
A
P





Lab Activity 1.1

1. Installing R

Latest version: 4.2.3 for windows

<https://cran.r-project.org/>

2. Installing R Studio

RStudio Desktop 2023.03.0+386

<https://posit.co/downloads/>

2: Install RStudio

[DOWNLOAD RSTUDIO DESKTOP FOR WINDOWS](#)

Size: 215.66 MB | [SHA-256: 93C7F307](#) | Version: 2023.12.0+369 |

Released: 2023-12-20



Assignment 1.1

1. Identify a DS project objective using the SMART methodology
2. Design your project using the SPAP methodology (Identify the DV and IVs)

Any
questions ?

vala.ali.rohani@estsetubal.ips.pt

