



# Introduction to Data Science

## Unit 1

# Course Contents

## Introduction to Data Science

[10 Hrs.]

Introduction to data science, Applications of data science; Limitations of data science Commonly used tools in data science, their strengths and common use-cases: R/RStudio, Python/Pandas/Jupyter Notebooks, Excel/Tableau/PowerBI;

Data Science life-cycle/Common methodologies for data science: CRISP-DM, OSEMN Framework, TDSP lifecycle;

Review of statistics and probability: Probability distributions, compound events and independence. Statistics: Centrality measures, variability measures, interpreting variance. Correlation analysis: Correlation coefficients, autocorrelation

# Before we begin.... Let's understand the data first

- Data is generated in various sources



Mobile Apps



Computer Applications



Point of Sale



Stock Market

Today	THU	FRI	SAT	SUN	MON	TUE
						
Sunny	Sunny	More sun than clouds	Passing clouds	More sun than clouds	Scattered clouds	Scattered clouds
66° 43°	69° 39°	72° 44°	78° 47°	78° 53°	77° 52°	75° 55°

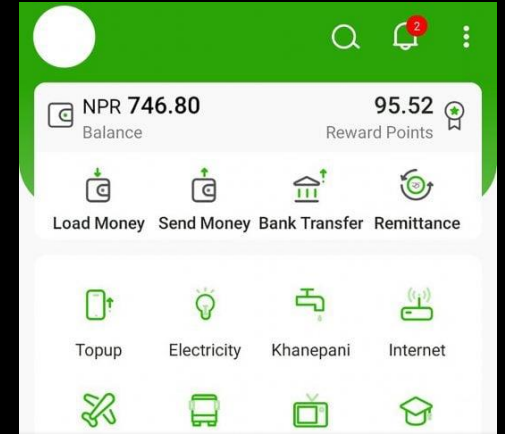
Weather



Security Surveillance

# Before we begin.... Let's understand the data first (contd.)

- Also,



## & many more



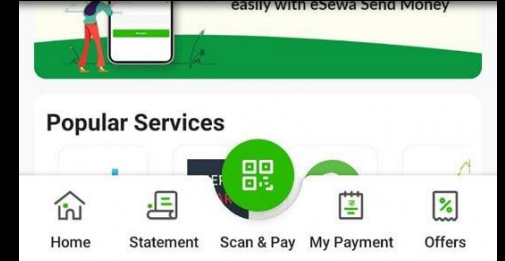
Astronomy



Medical Diagnosis



Scientific Study



Financial Transactions

# And, we categorize all these data sources into:

## INTERNAL DATA SOURCES:

Corporate ERP modules



Internal documents



Sensors, controllers



In-house call-centers



Website logs



## EXTERNAL DATA SOURCES:



Social media



Official statistics



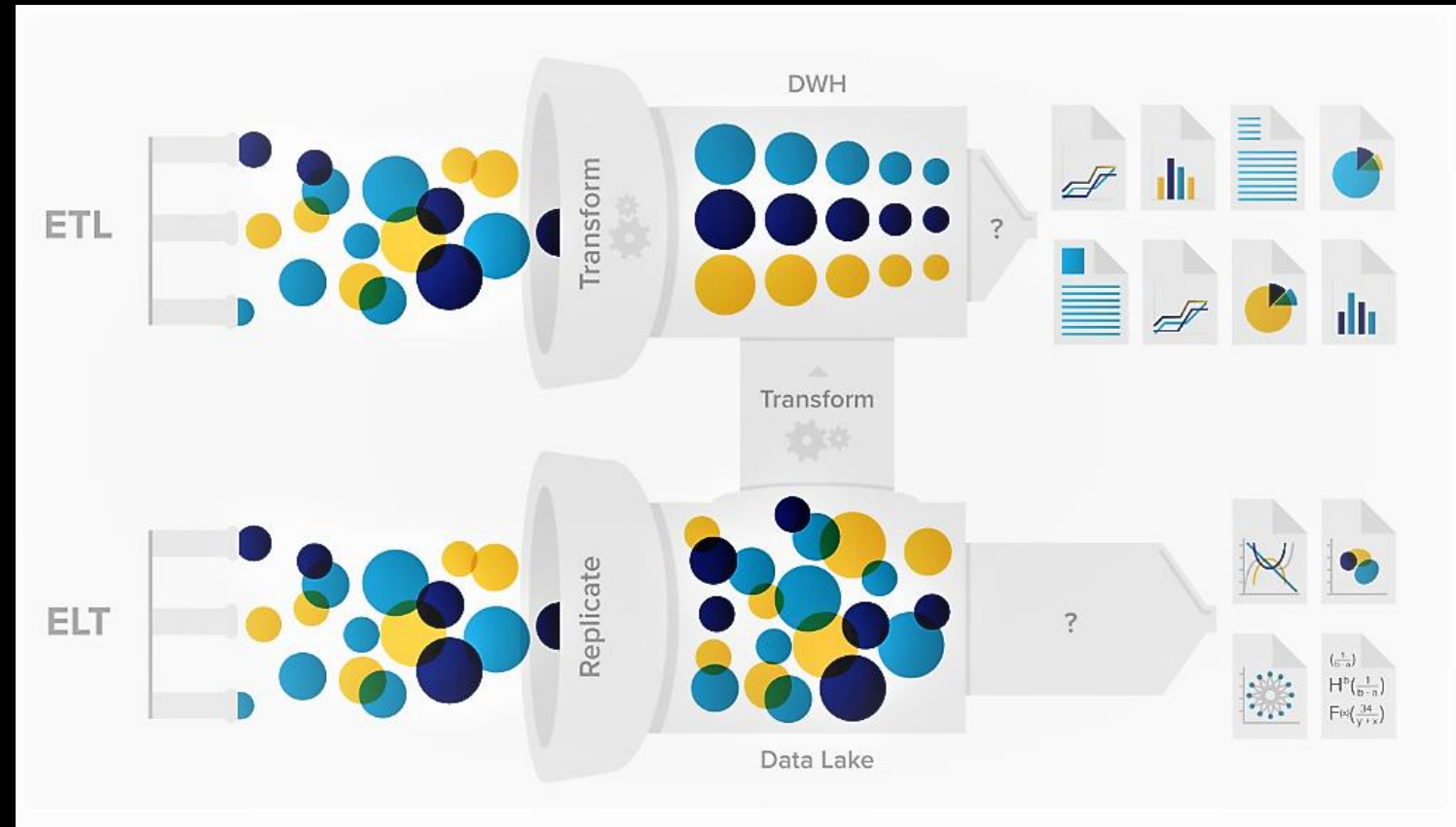
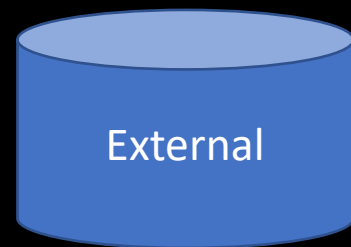
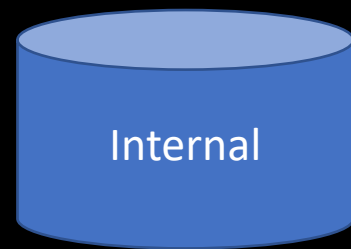
Weather forecasts



Publicly available data  
sets for machine learning



# Data migrates from various sources to Data Warehouse (Data Lake)



# Data Warehouse

- A data warehouse is a huge collection of business data used to help an organization make decisions.
- The large amount of data in data warehouses comes from different sources such as internal applications such as marketing, sales, and finance; customer-facing apps; and external partner systems, among others.
- On a technical level, a data warehouse periodically pulls data from those apps and systems; then, the data goes through formatting and import processes to match the data already in the warehouse.

# Data Warehouse (contd.)

- The data warehouse stores this processed data so it's ready for decision makers to access.
- How frequently data pulls occur, or how data is formatted, etc., will vary depending on the needs of the organization.

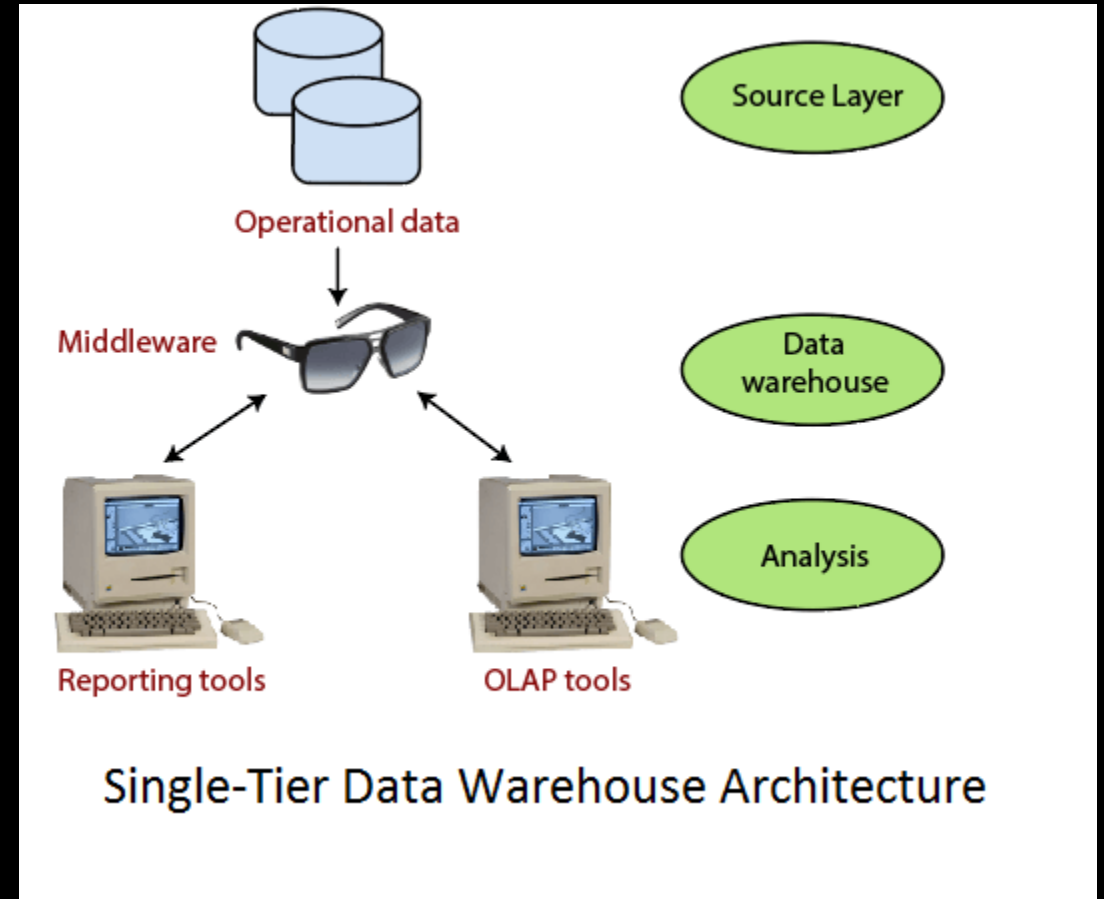


# Characteristics of Datawarehouse

1. **Subject Oriented**: focused on specific subject area
2. **Integrated**: Integrates data from multiple sources
3. **Time Variant**: Stores historical data
4. **Non Volatile**: Permanent Storage

# Data Warehouse Architecture

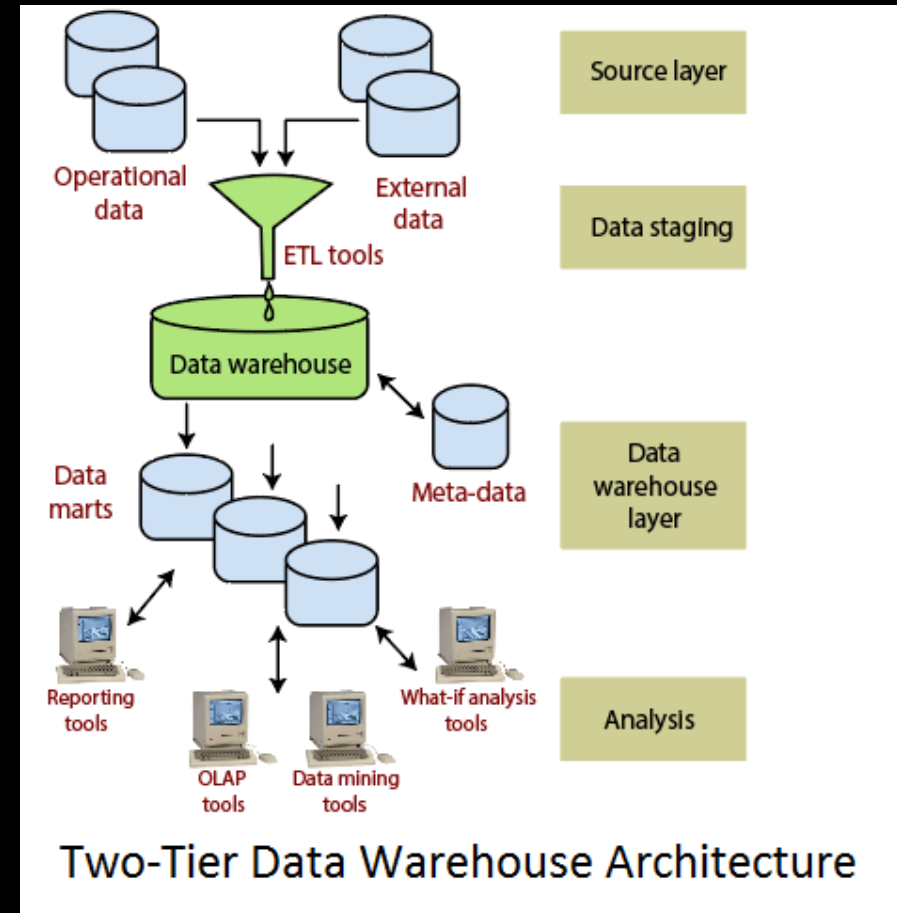
Single-Tier architecture is not periodically used in practice. Its purpose is to minimize the amount of data stored to reach this goal; it removes data redundancies.



# Data Warehouse Architecture (contd.)

Although it is typically called two-layer architecture to highlight a separation between physically available sources and data warehouses, in fact, consists of four subsequent data flow stages:

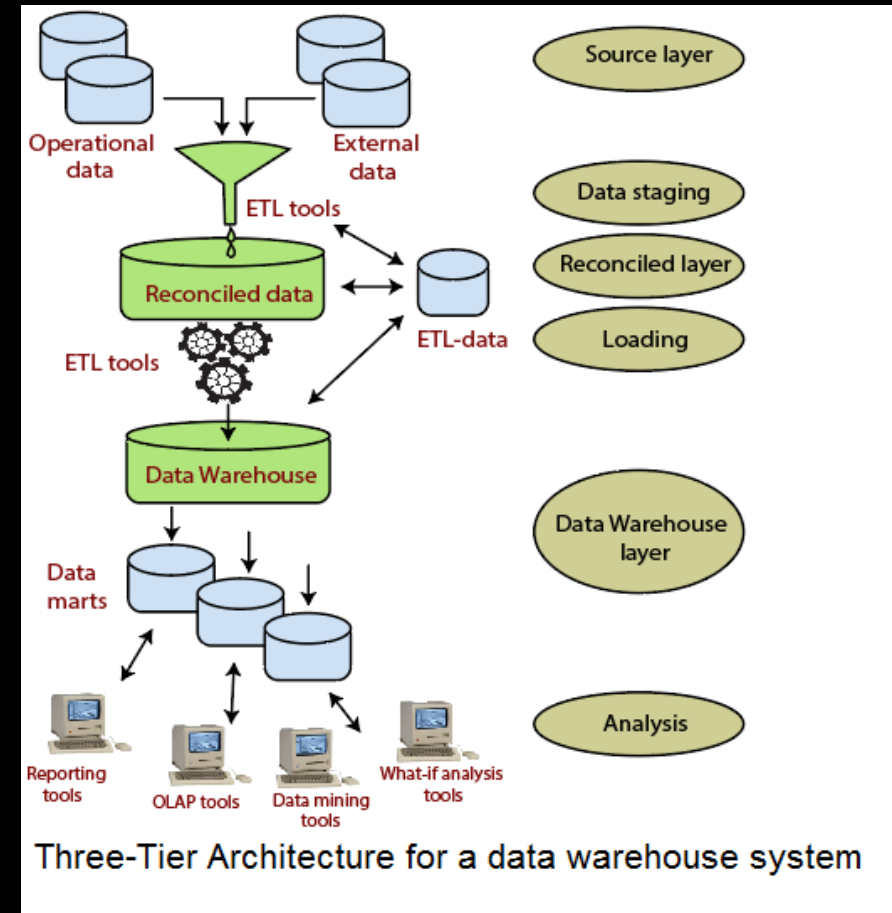
- a) Source layer
- b) Data Staging
- c) Data Warehouse layer
- d) Analysis



# Data Warehouse Architecture (contd.)

The three-tier architecture consists of the source layer (containing multiple source system), the reconciled layer and the data warehouse layer (containing both data warehouses and data marts). The reconciled layer sits between the source data and data warehouse.

The main advantage of the reconciled layer is that it creates a standard reference data model for a whole enterprise. At the same time, it separates the problems of source data extraction and integration from those of data warehouse population. In some cases, the reconciled layer is also directly used to accomplish better some operational tasks, such as producing daily reports that cannot be satisfactorily prepared using the corporate applications or generating data flows to feed external processes periodically to benefit from cleaning and integration.



# Data Lake

- A data lake is a central storage repository that holds big data from many resources in a raw, granular format.
- It can store structured, semi-structured or unstructured data, which means data can be kept in a more flexible format for future use.
- When storing data, a data lake associates it with identifiers and metadata tags for faster retrieval.
- The term “data lake” refers to the ad hoc nature of data in a data lake, as opposed to the clean and processed data stored in traditional data warehouse system.

# Data Lake (contd.)

- A data lake works on a principle called **schema-on-read**.
- This means that there is no predefined schema into which data needs to be fitted before storage.
- Only when the data is read during processing is it parsed and adapted into a schema as needed.
- This feature saves a lot of time that's usually spent on defining a schema. This also enables data to be stored as is, in any format.

*Data scientists can access, prepare, and analyze data faster and with more accuracy using data lakes. For analytics experts, this vast pool of data — available in various non-traditional formats — provides the opportunity to access the data for a variety of use cases like sentiment analysis or fraud detection.*



# Data lake vs data warehouse - Similarities

- A data lake and a data warehouse are similar in their basic purpose and objective, which make them easily confused:
  - Both are storage repositories that consolidate the various data stores in an organization.
  - The objective of both is to create a one-stop data store that will feed into various applications.

# Data lake vs data warehouse - Differences

- **Schema-on-read vs schema-on-write**

- The schema of a data warehouse is defined and structured before storage (schema is applied while writing data). A data lake, in contrast, has no predefined schema, which allows it to store data in its native format.
- In a data warehouse most of the data preparation usually happens before processing. In a data lake, it happens later, when the data is actually being used.

- **Complex vs simple user accessibility**

- As data is not organized in a simplified form before storage, a data lake often needs an expert with a thorough understanding of the various kinds of data and their relationships, to read through it. A data warehouse, in contrast, is easily accessible to both tech and non-tech users due its well-defined and documented schema. Even a new member on the team can begin to use a warehouse quickly.

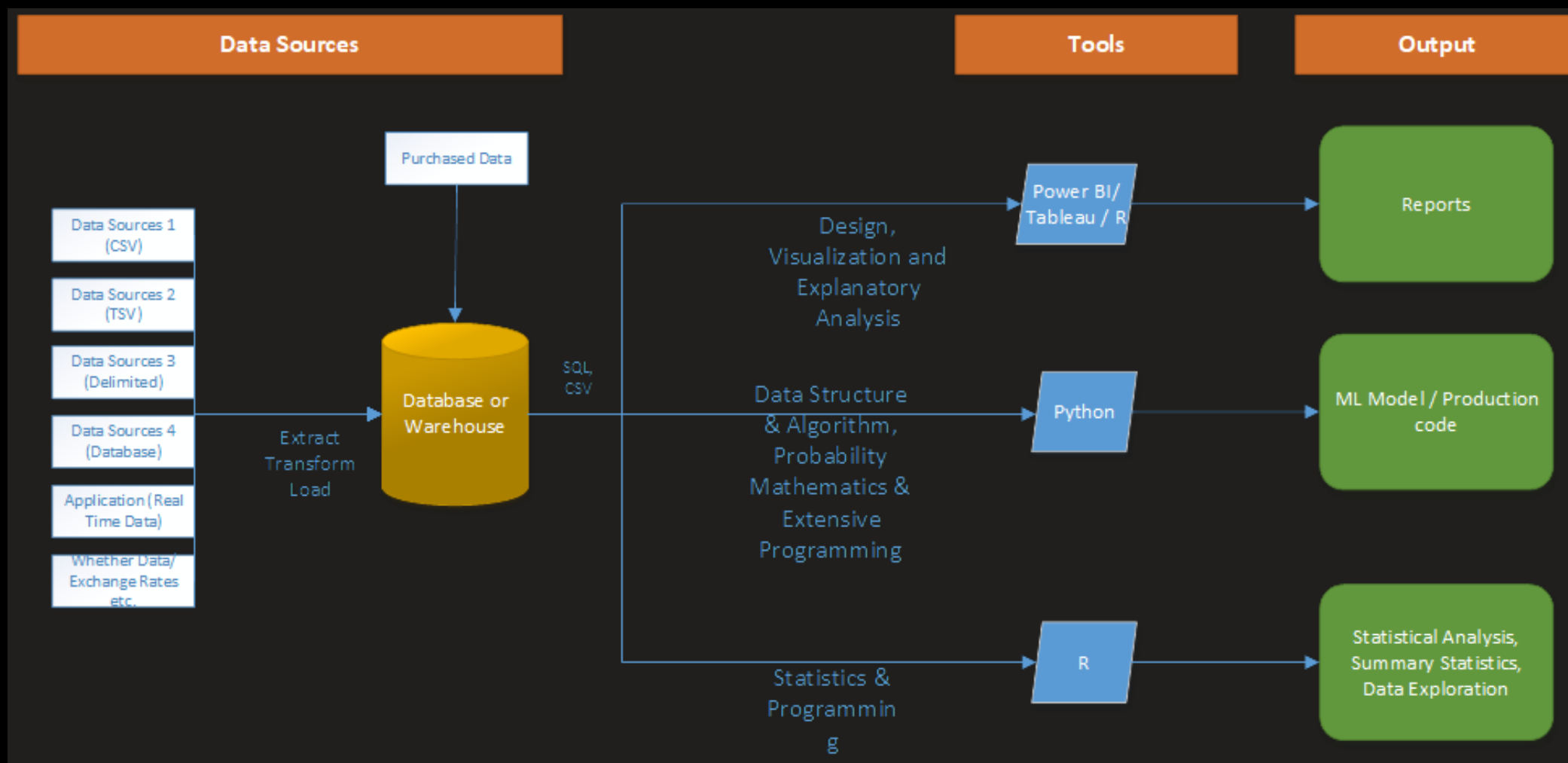
- **Flexibility vs rigidity**

- With a data warehouse, not only does it take time to define the schema at first, it also takes considerable resources to modify it when requirements change in the future. However, data lakes can adapt to changes easily. Also, as the need for storage capacity increases, it is easier to scale the servers on a data lake cluster.

# Data lake vs data warehouse - Differences

Characteristics	Data Warehouse	Data Lake
Data	Relational data from transactional systems, operational databases, and line of business applications	All data, including structured, semi-structured, and unstructured
Schema	Often designed prior to the data warehouse implementation but also can be written at the time of analysis  (schema-on-write or schema-on-read)	Written at the time of analysis (schema-on-read)
Price/Performance	Fastest query results using local storage	Query results getting faster using low-cost storage and decoupling of compute and storage
Data quality	Highly curated data that serves as the central version of the truth	Any data that may or may not be curated (i.e. raw data)
Users	Business analysts, data scientists, and data developers	Business analysts (using curated data), data scientists, data developers, data engineers, and data architects
Analytics	Batch reporting, BI, and visualizations	Machine learning, exploratory analytics, data discovery, streaming, operational analytics, big data, and profiling

# Thus,



# Introduction to Data Science

- Over the past few years, there's been a lot of hype in the internet about "data science" and "Big Data."
- But, what actually "Data Science" is? And what does "Big Data" means?
- How "Data Science" and "Big Data" are related?
- Is data science the science of Big Data?
- Is data science only the stuff going on in companies like Google and Facebook and tech companies?
- Why do many people refer to Big Data as crossing disciplines (astronomy, finance, tech, etc.) and to data science as only taking place in tech? Just how big is big? Or is it just a relative term?

# Introduction to Data Science (contd.)

## Is it new Age Thing?

- Statisticians already feel that they are studying and working on the “Science of Data.” That’s their bread and butter.
- Many of the algorithms ( *e.g. linear regression, logistic regression, Bayesian Statistics, and even neural network* ) we used today were discovered long back in the past.
- However, many of the methods and techniques we’re using—and the challenges we’re facing now—are part of the evolution of everything that’s come before.



# Introduction to Data Science (contd.)

## Why now?

- We have massive amounts of data about many aspects of our lives, and, simultaneously, an abundance of inexpensive computing power.
- *Shopping, communicating, reading news, listening to music, searching for information, expressing our opinions*—all this is **being tracked online**, as most people know.
- What people might not know is that the **“datafication”** of our **offline behavior** has started as well, mirroring the online data collection revolution.
- Put the two together, and there’s a lot to learn about our behavior and, by extension, who we are as a species.

# Introduction to Data Science (contd.)

## Why now?

- It's not just Internet data, though—it's finance, the medical industry, pharmaceuticals, bioinformatics, social welfare, government, education, retail, and the list goes on.
- But it's not only the massiveness that makes all this new data interesting (or poses challenges).
- It's that the data itself, often in real time, becomes the building blocks of data products.

# Introduction to Data Science (contd.)

## Why now?

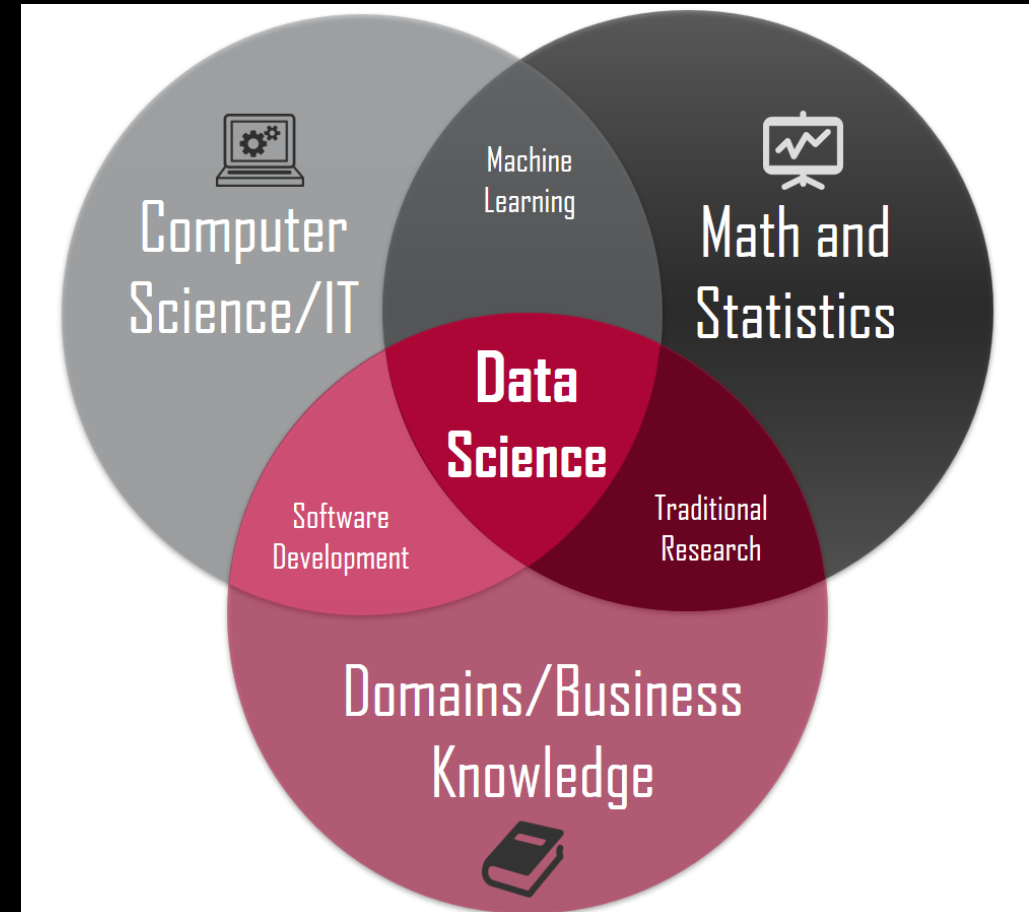
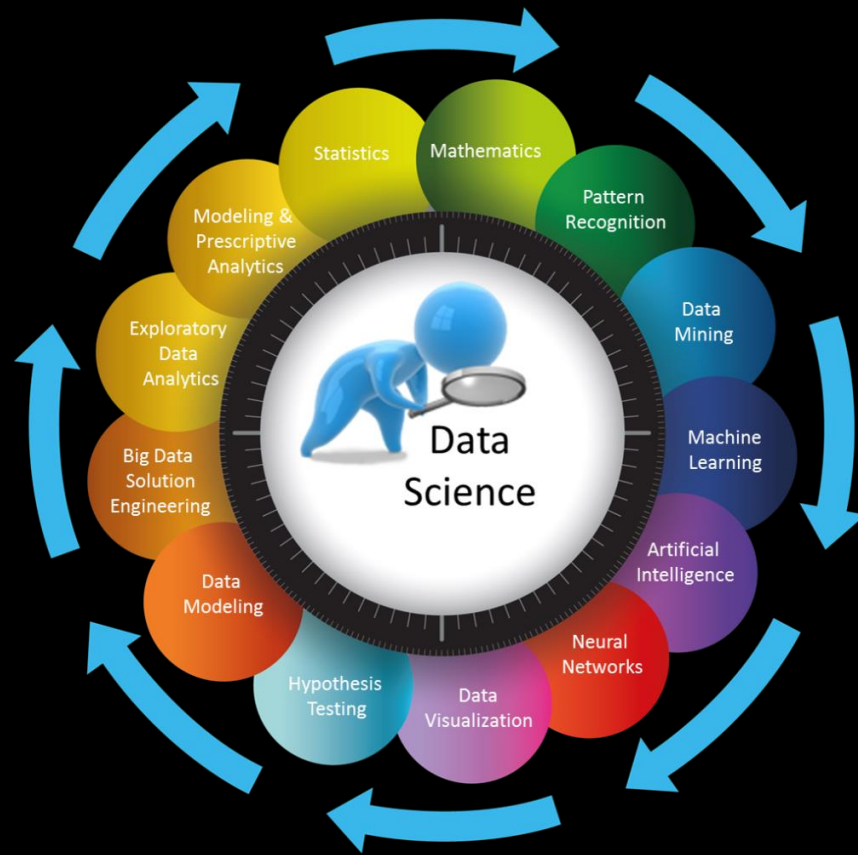
- We're witnessing the beginning of a massive, culturally saturated feedback loop where our behavior changes the product and the product changes our behavior.
- Technology makes this possible: infrastructure for large-scale data processing, increased memory, and bandwidth, as well as a cultural acceptance of technology in the fabric of our lives.

# Introduction to Data Science (contd.)

## So, what is data science?

- Is it new, or is it just statistics or analytics rebranded? Is it real, or is it pure hype? And if it's new and if it's real, what does that mean?
- This is an ongoing discussion, and many people have their own definition, their own explanations.

# Introduction to Data Science (contd.)



*Drew Conway's Venn diagram of data science*

# Introduction to Data Science (contd.)

- Metamarket CEO Mike Driscoll defines data science as:
  - Data science, as it's practiced, is a blend of Red-Bull-fueled hacking and espresso-inspired statistics.
  - But data science is not merely hacking—because when hackers finish debugging their Bash one-liners and Pig scripts, few of them care about non-Euclidean distance metrics.
  - And data science is not merely statistics, because when statisticians finish theorizing the perfect model, few could read a tab-delimited file into R if their job depended on it.
  - Data science is the civil engineering of data. Its acolytes possess a practical knowledge of tools and materials, coupled with a theoretical understanding of what's possible.



# Introduction to Data Science (contd.)

- Who is data scientist?

A data scientist is someone who knows how to extract meaning from and interpret data, which requires both tools and methods from statistics and machine learning, as well as being human. She spends a lot of time in the process of collecting, cleaning, and munging data, because data is never clean. This process requires persistence, statistics, and software engineering skills—skills that are also necessary for understanding biases in the data, and for debugging logging output from code.

# Introduction to Data Science (contd.)

- Who is data scientist?

A **chief data scientist** should be **setting the data strategy of the company**, which involves a variety of things: setting everything up from the engineering and infrastructure for collecting data and logging, to privacy concerns, to deciding what data will be user-facing, how data is going to be used to make decisions, and how it's going to be built back into the product. She should manage a team of engineers, scientists, and analysts and should communicate with leadership across the company, including the CEO, CTO, and product leadership. She'll also be concerned with patenting innovative solutions and setting research goals.

# Introduction to Data Science (contd.)

- Who is data scientist?

Once she gets the data into shape, a crucial part is exploratory data analysis, which combines visualization and data sense. She'll find patterns, build models, and algorithms—some with the intention of understanding product usage and the overall health of the product, and others to serve as prototypes that ultimately get baked back into the product. She may design experiments, and she is a critical part of data driven decision making. She'll communicate with team members, engineers, and leadership in clear language and with data visualizations so that even if her colleagues are not immersed in the data themselves, they will understand the implications.

# Applications of Data Science

## 1. Finance

- Risk Management
- Financial Analysis
- Stock price Prediction

## 2. Retailing

- Inventory Replenishment
- Customer Segmentation
- Promotion campaigning
- Sales Campaigning
- Demand and Supply projection

## 3. Banking

- Resource Utilization
- Fraud Detection
- Loan Management
- Predictive Analysis
- Customer Segmentation

## 4. Manufacturing

- Resource Utilization
- Monitoring of energy costs and optimize production hour.
- Identify potential problems through analysis of continuous stream of data

# Applications of Data Science (contd.)

## 5. Health

- Medical Image Analysis
- Genetics and Genomics
- Drug Discovery
- Predictive Modeling for Diagnosis
- Health bots or virtual assistants

## • Other Applications

- Fraud and Risk Detection
- Personalized Advertising
- Content Recommendation
- Advanced Image Recognition
- Airline Route Planning
- Gaming etc.

# Limitations of Data Science

- It's a Blurry Team
- Mastering Data Science is near to impossible.
- Large amount of Domain knowledge is required
- Arbitrary data may yield unexpected result.
- Data alone cannot guide decisions because they are a means of orienting you toward
- Difficulty in devising a conclusion if dataset is smaller.
- Garbage in Garbage Out



# Commonly used tools

- Python
- R
- Tableau
- Power BI
- Excel / Google Sheets
- Jupyter
- Weka
- Apache Spark
- Apache Hadoop
- MATLAB
- Julia
- Scala
- Azure/Google Cloud/AWS
- MLOps
- AutoML

# Commonly used library

## Python

- NumPy
- SciPy
- Pandas
- Matplotlib
- SciKit-Learn
- Tensorflow
- Pytorch
- BeautifulSoup
- NLTK
- OpenCV
- ScraPY
- XGBoost
- Seaborn

## R

- Dplyr
- TidyR
- Readr
- Stringr
- Ggplot2
- Lubridate
- Jsonlite
- BioConductor
- Shiny
- Knitr
- Caret
- Rmarkdown

# Basic Steps for Data Science



Before we begin...

# Understand the convention

Given is our dataset  $D$

$x_1$	$x_2$	...	$x_n$	$y$

- $(x^{(i)}, y^{(i)})$  is the  $i$ -th sample from our dataset
- $x_j$  is the  $j$ -th feature from our dataset
- $X$  is the  $m * n$  matrix for collections of features
- $y$  is the  $m * 1$  vector for target ( labels )

# Let's make hands dirty

Task 01 – Build Classification Model to classify Survival in Titanic

URL - [shorturl.at/grxBJ](https://shorturl.at/grxBJ)

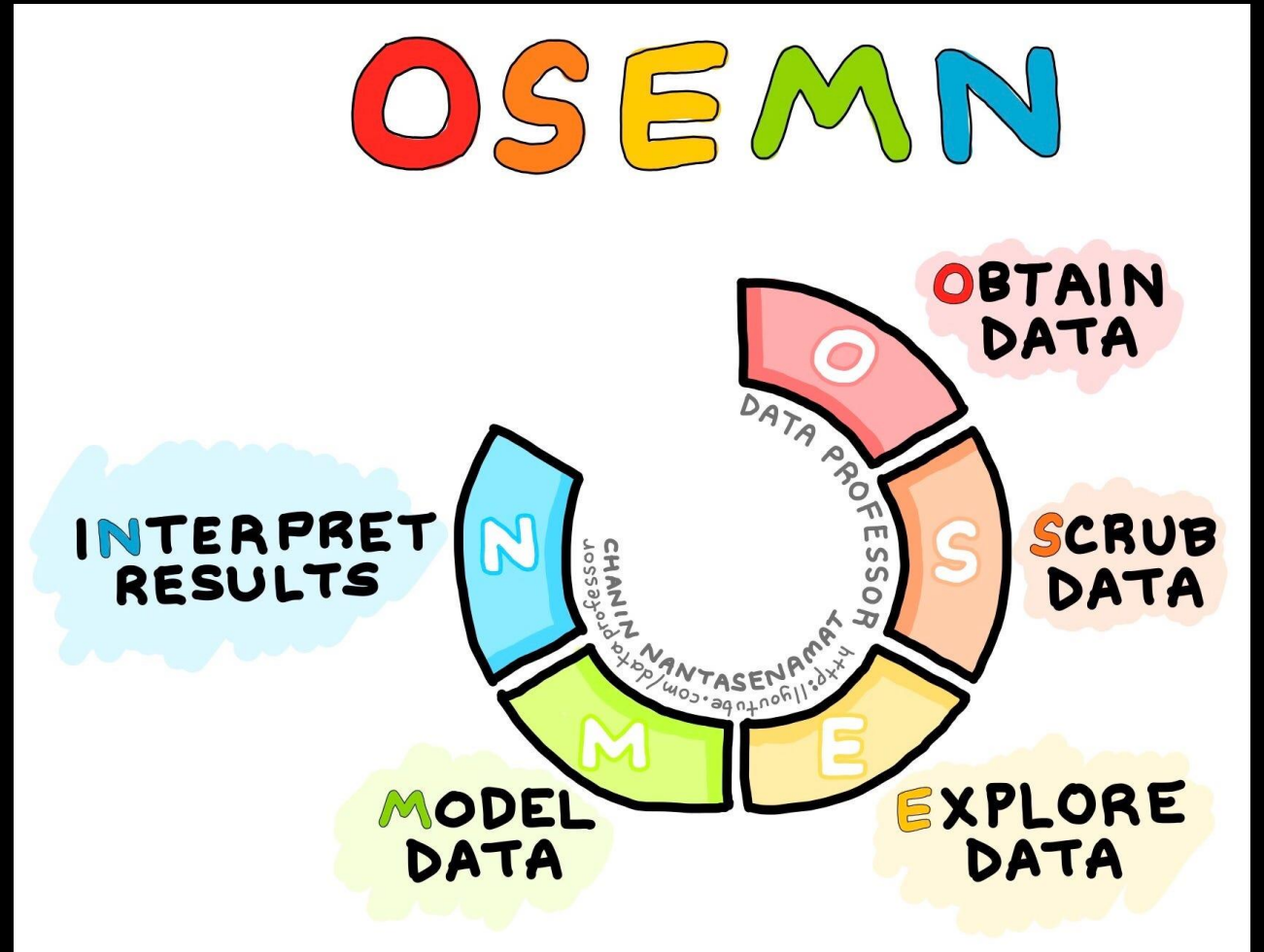
# Data Science Life Cycle

# OSEMN Framework



# OSEMN Framework

- Obtain
- Scrub
- Explore
- Model
- iNterpret



# OSEMN Framework (Contd.)

- Benefits:
  - **Simple** – It distills the complex process of a data science project into five clear steps. This is especially noteworthy given that this general process and the modern concept of data science were still new when Mason and Wiggins created OSEMN in 2010.
  - **Catchy** – OSEMN is Awesome!
  - **Makes sense** – The steps presented have a logical flow representative of the general data science life cycle.
  - **Provides a shared understanding** – OSEMN creates a taxonomy to help define how a data science project progresses.

# OSEMN Framework (Contd.)

- Shortcomings:
  - **Misses business understanding** – The framework starts with Obtain which ignores the key base questions that should come first, namely: “Should I invest time on this project?” and “What outcome am I trying to drive?”
  - **Doesn’t consider deployment** – OSEMN implicitly assumes that you are delivering a one-time output. In reality, you often need to deploy a model in a production system so that it continues to provide value over time.
  - **Ignores teamwork** – Data science is increasingly a team sport. Yet, OSEMN ignores the broader team aspect of modern projects.
  - **It’s linear** – OSEMN proceeds in a waterfall-like manner with each phase following the other. In reality, you often switch back and forth between phases as needed. Moreover, you will want frequent decision points where you re-assess and adjust your plan based on recent learnings.

# Team Data Science Process (TDSP)

# TDSP

- The Team Data Science Process (TDSP) is an agile, iterative data science methodology to deliver predictive analytics solutions and intelligent applications efficiently.
- TDSP helps improve team collaboration and learning by suggesting how team roles work best together.
- TDSP includes best practices and structures from Microsoft and other industry leaders to help toward successful implementation of data science initiatives.
- The goal is to help companies fully realize the benefits of their analytics program.

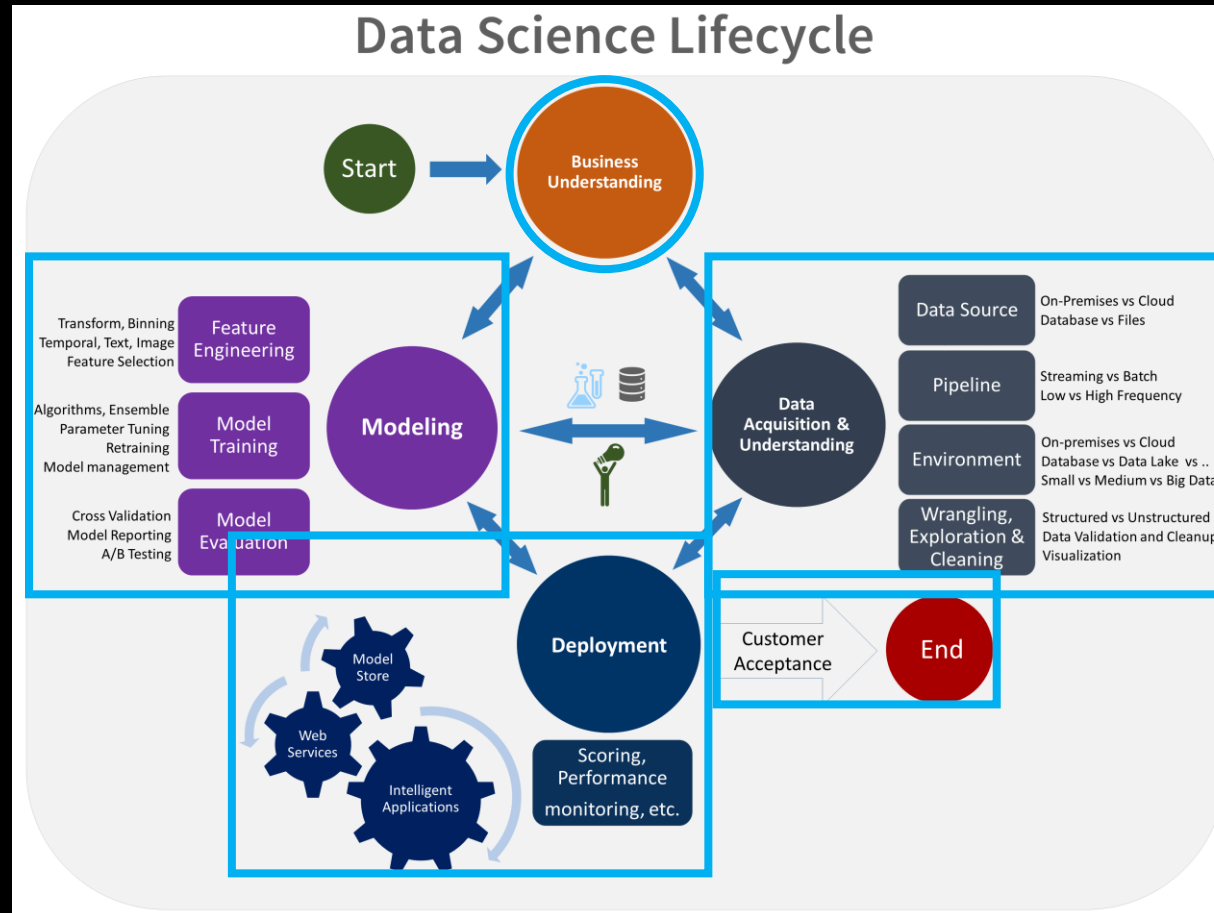
# TDSP (contd.)

- Key Components of TDSP
  - A data science lifecycle
  - A standardized project structure
  - Infrastructure and resources
  - Tools and utilities

# TDSP (contd.)

- Data Science Lifecycle
  - The Team Data Science Process (TDSP) provides a lifecycle to structure the development of your data science projects.
  - The lifecycle outlines the full steps that successful projects follow.
  - The lifecycle outlines the major stages that projects typically execute, often iteratively:
    - Business Understanding
    - Data Acquisition and Understanding
    - Modeling
    - Deployment

# TDSP (contd.)

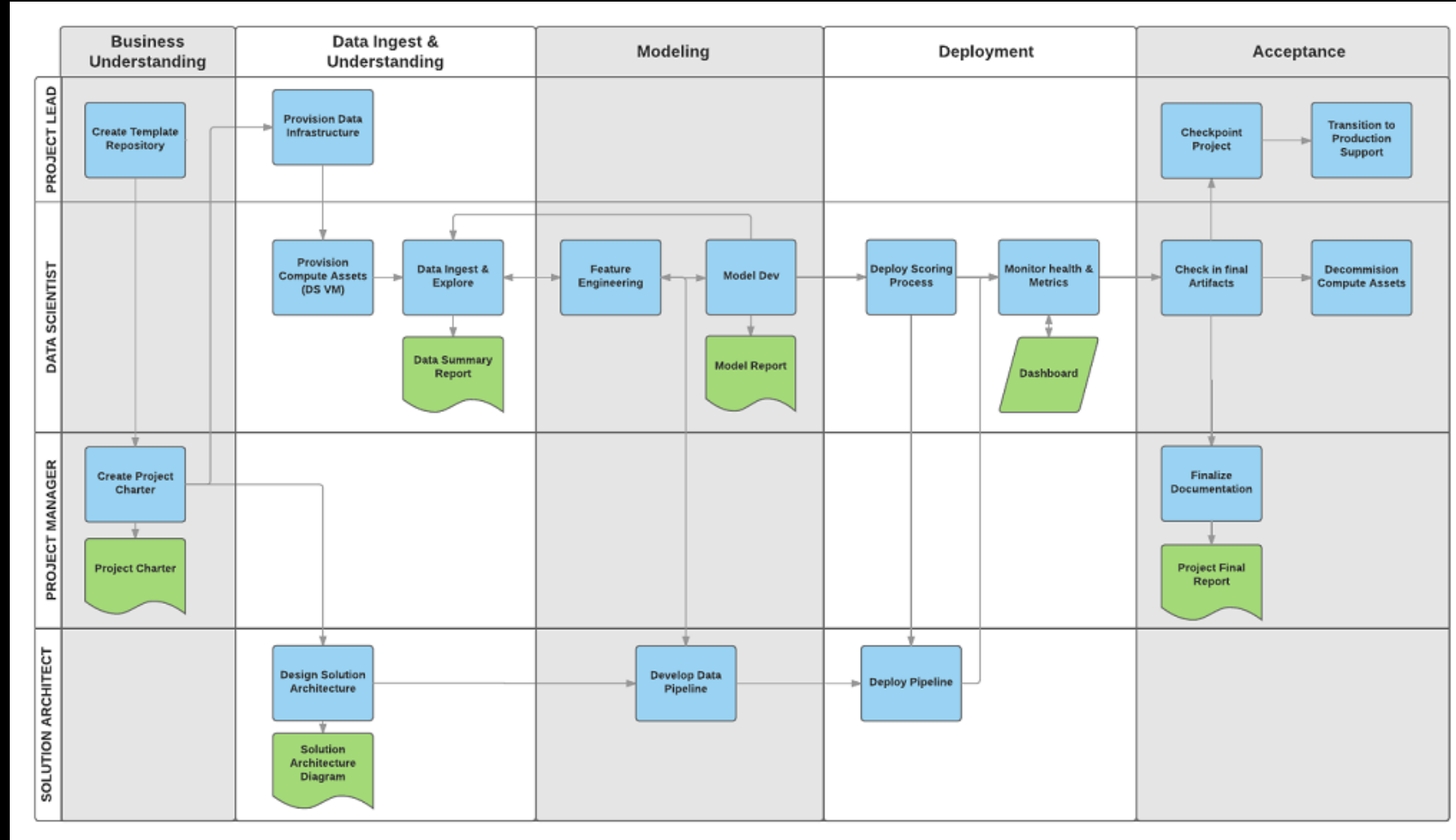




# TDSP (contd.)

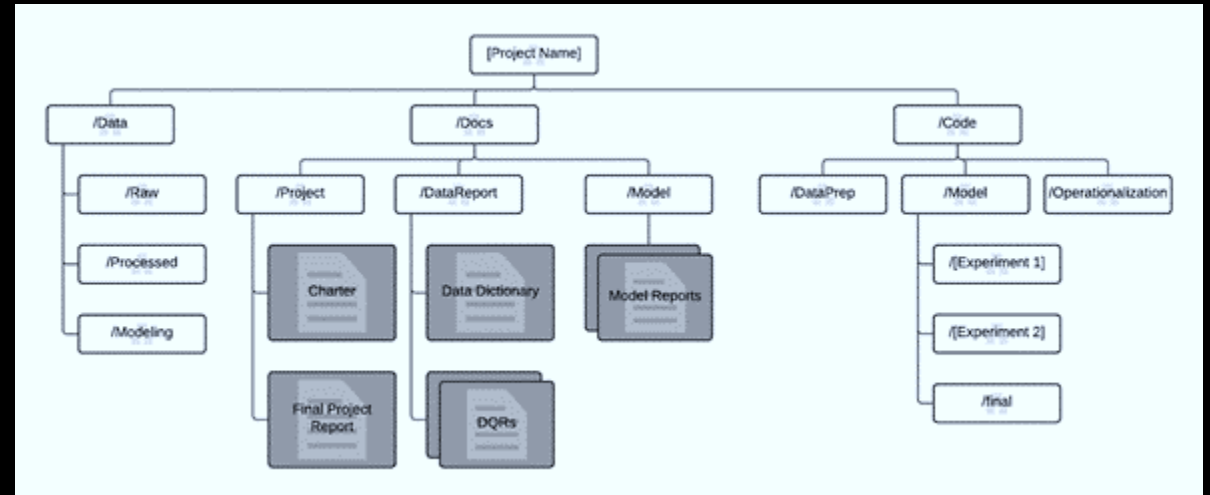
- The goals, tasks, and documentation artifacts for each stage of the lifecycle in TDSP are described in the Team Data Science Process lifecycle topic.
- These tasks and artifacts are associated with project roles:
  - Solution architect
  - Project manager
  - Data engineer
  - Data scientist
  - Application developer
  - Project lead

# TDSP (contd.)



# TDSP (contd.)

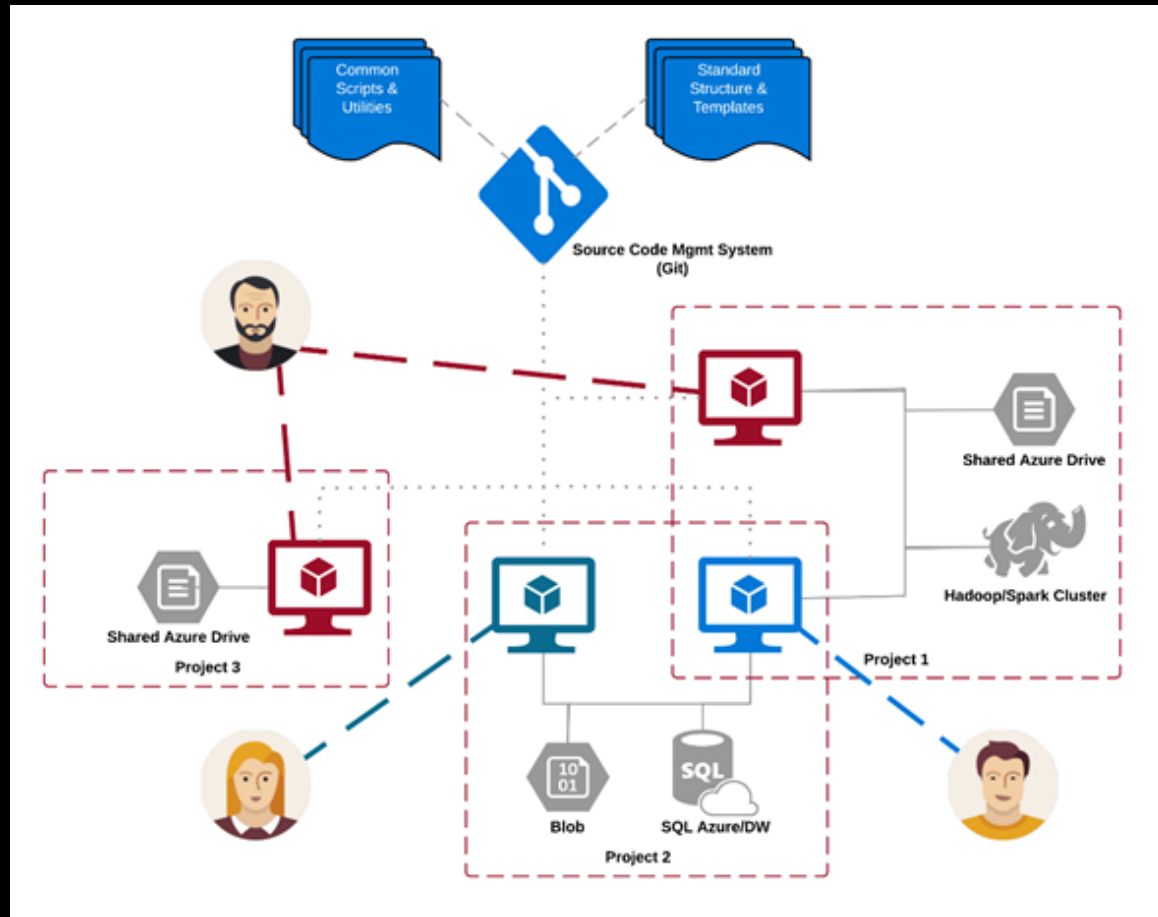
- **Standardized project structure**
  - Having all projects share a directory structure and use templates for project documents makes it easy for the team members to find information about their projects.
  - All code and documents are stored in a version control system (VCS) like Git, TFS, or Subversion to enable team collaboration.
  - Tracking tasks and features in an agile project tracking system like Jira, Rally, and Azure DevOps allows closer tracking of the code for individual features. Such tracking also enables teams to obtain better cost estimates.
  - TDSP recommends creating a separate repository for each project on the VCS for versioning, information security, and collaboration. The standardized structure for all projects helps build institutional knowledge across the organization.



# TDSP (contd.)

- Infrastructure and resources for data science projects
  - TDSP provides recommendations for managing shared analytics and storage infrastructure such as:
    - cloud file systems for storing datasets
    - databases
    - big data (SQL or Spark) clusters
    - machine learning service
  - The analytics and storage infrastructure, where raw and processed datasets are stored, may be in the cloud or on-premises.
  - This infrastructure enables reproducible analysis.
  - It also avoids duplication, which may lead to inconsistencies and unnecessary infrastructure costs. Tools are provided to provision the shared resources, track them, and allow each team member to connect to those resources securely.
  - It is also a good practice to have project members create a consistent compute environment. Different team members can then replicate and validate experiments.

# TDSP (contd.)



Here is an example of a team working on multiple projects and sharing various cloud analytics infrastructure components.

# TDSP (contd.)

- Tools and utilities for project execution
  - Introducing processes in most organizations is challenging.
  - Tools provided to implement the data science process and lifecycle help lower the barriers to and increase the consistency of their adoption.
  - TDSP provides an initial set of tools and scripts to jump-start adoption of TDSP within a team.
  - It also helps automate some of the common tasks in the data science lifecycle such as data exploration and baseline modeling.
  - There is a well-defined structure provided for individuals to contribute shared tools and utilities into their team's shared code repository.
  - These resources can then be leveraged by other projects within the team or the organization.
  - Microsoft provides extensive tooling inside Azure Machine Learning supporting both open-source (Python, R, ONNX, and common deep-learning frameworks) and also Microsoft's own tooling (AutoML).

# TDSP (contd.)

## Benefits

- Elaborated Documentation
- Agile
- Familiar
- Data Science Native
- Flexible
- Detailed
- Free Templates

## Short comings

- Fixed Sprints: TDSP leverages fixed-length planning sprints which many data scientists struggle with.
- Some Inconsistencies: Not all of Microsoft's documentation is consistent.

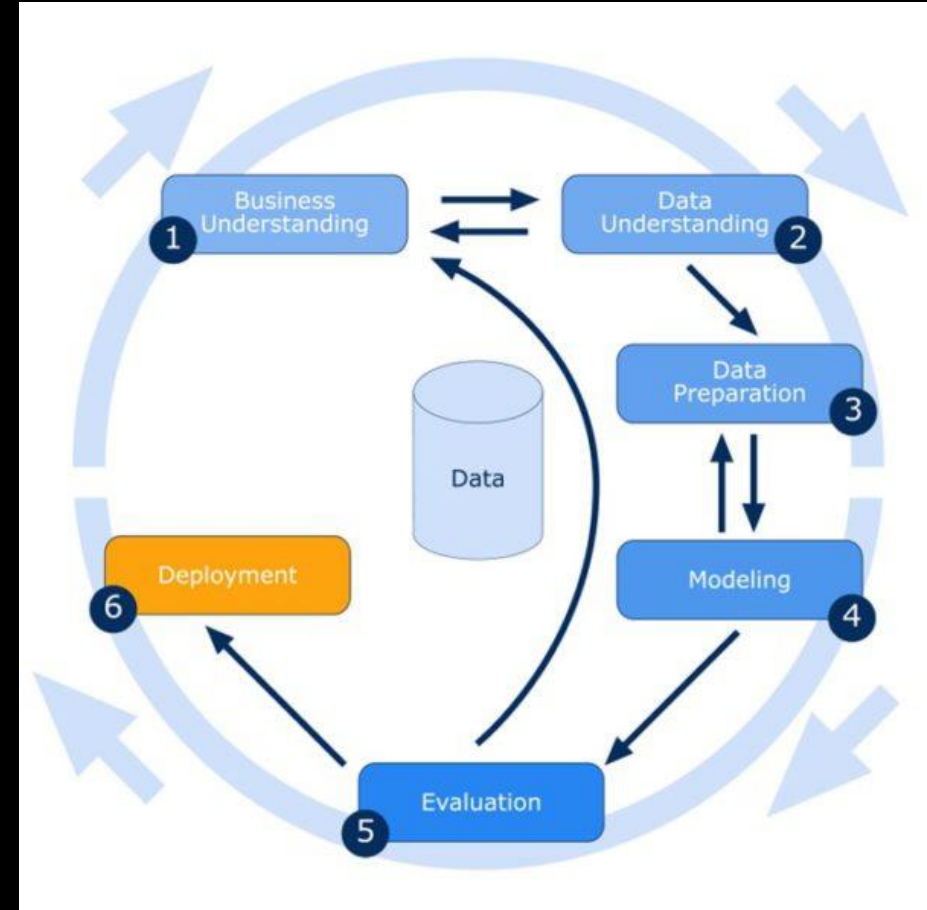
# TDSP (contd.)





# CRISP-DM

- The **C**ross **I**ndustry **S**tandard **P**rocess for **D**ata **M**ining (CRISP-DM) is a process model that serves as the base for a data science process. It has six sequential phases:
  - Business understanding
  - Data understanding
  - Data preparation
  - Modeling
  - Evaluation
  - Deployment



<https://www.datascience-pm.com/crisp-dm-2/>

<https://web.archive.org/web/20220401041957/https://www.the-modeling-agency.com/crisp-dm.pdf>

# CRISP-DM (contd.)

## I. Business Understanding

- Any good project starts with a deep understanding of the customer's needs. Data mining projects are no exception and CRISP-DM recognizes this.
- The Business Understanding phase focuses on understanding the objectives and requirements of the project. Aside from the third task, the three other tasks in this phase are foundational project management activities that are universal to most projects:
  1. **Determine business objectives:** You should first “thoroughly understand, from a business perspective, what the customer really wants to accomplish.” (CRISP-DM Guide) and then define business success criteria.
  2. **Assess situation:** Determine resources availability, project requirements, assess risks and contingencies, and conduct a cost-benefit analysis.
  3. **Determine data mining goals:** In addition to defining the business objectives, you should also define what success looks like from a technical data mining perspective.
  4. **Produce project plan:** Select technologies and tools and define detailed plans for each project phase.

# CRISP-DM (contd.)

## II. Data Understanding

- Next is the Data Understanding phase. Adding to the foundation of Business Understanding, it drives the focus to identify, collect, and analyze the data sets that can help you accomplish the project goals. This phase also has four tasks:
  1. **Collect initial data**: Acquire the necessary data and (if necessary) load it into your analysis tool.
  2. **Describe data**: Examine the data and document its surface properties like data format, number of records, or field identities.
  3. **Explore data**: Dig deeper into the data. Query it, visualize it, and identify relationships among the data.
  4. **Verify data quality**: How clean/dirty is the data? Document any quality issues.

# CRISP-DM (contd.)

## III. Data Preparation

- A common rule of thumb is that 80% of the project is data preparation.
- This phase, which is often referred to as “data munging”, prepares the final data set(s) for modeling. It has five tasks:
  1. **Select data**: Determine which data sets will be used and document reasons for inclusion/exclusion.
  2. **Clean data**: Often this is the lengthiest task. Without it, you’ll likely fall victim to garbage-in, garbage-out. A common practice during this task is to correct, impute, or remove erroneous values.
  3. **Construct data**: Derive new attributes that will be helpful. For example, derive someone’s body mass index from height and weight fields.
  4. **Integrate data**: Create new data sets by combining data from multiple sources.
  5. **Format data**: Re-format data as necessary. For example, you might convert string values that store numbers to numeric values so that you can perform mathematical operations.

# CRISP-DM (contd.)

## IV. Modeling

- What is widely regarded as data science's most exciting work is also often the shortest phase of the project.
- Here you'll likely build and assess various models based on several different modeling techniques. This phase has four tasks:
  1. **Select modeling techniques**: Determine which algorithms to try (e.g. regression, neural net).
  2. **Generate test design**: Pending your modeling approach, you might need to split the data into training, test, and validation sets.
  3. **Build model**: As glamorous as this might sound, this might just be executing a few lines of code like `reg = LinearRegression().fit(X, y)`.
  4. **Assess model**: Generally, multiple models are competing against each other, and the data scientist needs to interpret the model results based on domain knowledge, the pre-defined success criteria, and the test design.

# CRISP-DM (contd.)

## V. Evaluation

- Whereas the Assess Model task of the Modeling phase focuses on technical model assessment, the Evaluation phase looks more broadly at which model best meets the business and what to do next. This phase has three tasks:
  1. **Evaluate results:** Do the models meet the business success criteria? Which one(s) should we approve for the business?
  2. **Review process:** Review the work accomplished. Was anything overlooked? Were all steps properly executed? Summarize findings and correct anything if needed.
  3. **Determine next steps:** Based on the previous three tasks, determine whether to proceed to deployment, iterate further, or initiate new projects.

# CRISP-DM (contd.)

## VI. Deployment

*“Depending on the requirements, the deployment phase can be as simple as generating a report or as complex as implementing a repeatable data mining process across the enterprise.”*

— CRISP-DM Guide

- A model is not particularly useful unless the customer can access its results. The complexity of this phase varies widely. This final phase has four tasks:
  - **Plan deployment**: Develop and document a plan for deploying the model.
  - **Plan monitoring and maintenance**: Develop a thorough monitoring and maintenance plan to avoid issues during the operational phase (or post-project phase) of a model.
  - **Produce final report**: The project team documents a summary of the project which might include a final presentation of data mining results.
  - **Review project**: Conduct a project retrospective about what went well, what could have been better, and how to improve in the future.

# Is CRISP-DM Agile or Waterfall?

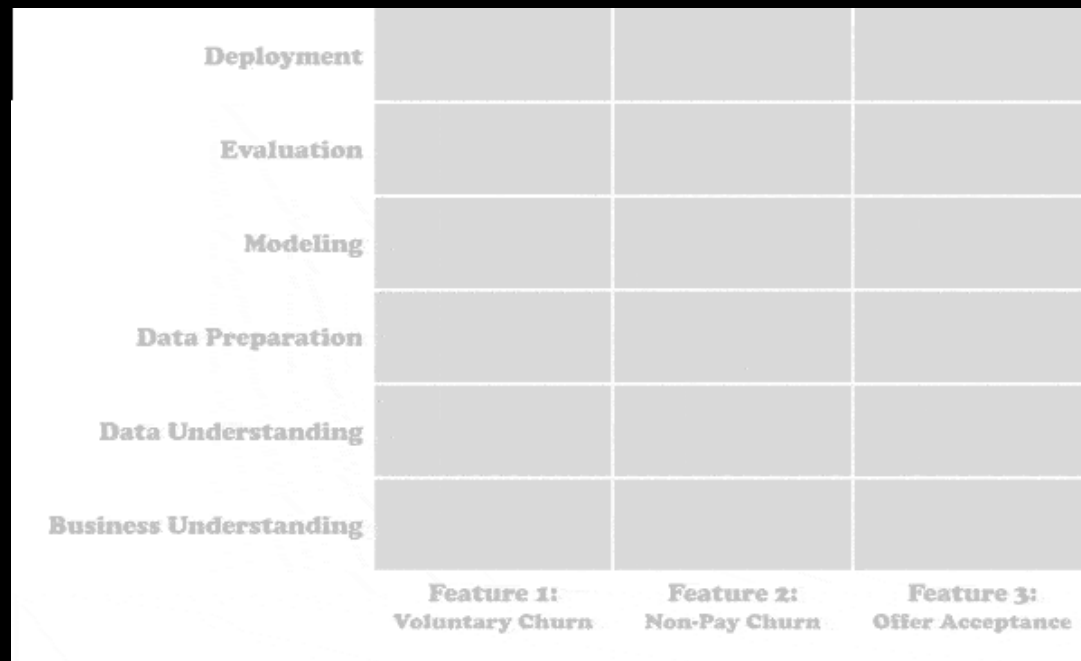
- Some argue that it is flexible and agile and while others see CRISP-DM as rigid. What really matters is how you implement it.
- **Waterfall**
  - On one hand, many view CRISP-DM as a rigid waterfall process – in part because of its reporting requirements are excessive for most projects. Moreover, the guide states in the business understanding phase that “the project plan contains detailed plans for each phase” – a hallmark aspect of traditional waterfall approaches that require detailed, upfront planning.
  - Indeed, if you follow CRISP-DM precisely (defining detailed plans for each phase at the project start and include every report) and choose not to iterate frequently, then you’re operating more of a waterfall process.
- **Agile**
  - On the other hand, CRISP-DM indirectly advocates agile principles and practices by stating: “The sequence of the phases is not rigid. Moving back and forth between different phases is always required. The outcome of each phase determines which phase, or particular task of a phase, has to be performed next.”
  - Thus if you follow CRISP-DM in a more flexible way, iterate quickly, and layer in other agile processes, you’ll wind up with an agile approach.

*Example: To illustrate how CRISP-DM could be implemented in either an Agile or waterfall manner, imagine a churn project with three deliverables: a voluntary churn model, a non-pay disconnect churn model, and a propensity to accept a retention-focused offer.*



# CRISP-DM (contd.)

## Waterfall: Horizontal Slicing



## Agile: Vertical Slicing



# CRISP-DM (contd.)

## Which is better?

- When possible, take an agile approach and slice vertically so that:
  - Stakeholders get value sooner
  - Stakeholders can provide meaningful feedback
  - The data scientists can assess model performance earlier
  - The project team can adjust the plan based on stakeholder feedback

# CRISP-DM (contd.)

## Benefits

- Generalizable
- Common Sense
- Adoptable
- Right Start
- Strong Finish
- Flexible

## Shortcomings

- Rigid
- Documentation Heavy
- Not Modern
- Not a Project Management Approach