

Data Science with Python

Day 1 - Introduction to Data Science

Today's Agenda

✓ Data Science Introduction

- What is Data Science?
- Some Examples of Data Science Project Objective

✓ Data Science Toolkit

- What you need to learn?

✓ Job outlook

- What are the opportunities?

✓ Target Audience & Prerequisites

- Who is more suitable?
- What is expected?

✓ Your Specific Questions

Data Science Introduction

What is Data Science?

“Data science, also known as data-driven science, is an interdisciplinary field about scientific methods, processes, and systems to extract knowledge or insights from data present in various forms, either structured or unstructured.”

Data science employs **techniques** and **theories** drawn from a wide range of disciplines like **Mathematics**, **Statistics**, **Information Science**, and **Computer Science**, in particular from the subdomains of **Machine learning**, **Classification**, **Association**, **Cluster analysis**, **Data mining**, **forecasting** and **Visualization** in order to understand and analyze actual phenomena with data.

Data Science : Data + Science = Knowledge

Examples of Data Science Project Objectives

- **Regression Analysis** – Finding the relationship between a dependent variable and one or more independent variables - Predicting Diamond price based on Carat, Cut & Clarity
- **Classification Analysis** – Dividing objects into 2 or more known classes - Distinguishing cancer and normal cells
- **Anomalies Detection(Outliers Analysis)** - Finding unusual - Credit card transactions
- **Association Analysis** – Finding links - Shopping cart analysis
- **Cluster Analysis (Segmentation)** – Grouping similar objects together - Grouping customers into different clusters based on their previous shopping data/transactions.
- **Time Series Analysis** – Time dependent Data - Stock prediction

Data Science Toolkit

What you need to learn?

- ✓ **Data Engineering** – Getting and Processing the data - Python, R, PL/SQL, SAS, ETL, Big Data
- ✓ **Data Analysis** – Exploratory Data Analysis, Predictive and Descriptive Analysis - Python, R, SAS
- ✓ **Statistics & Probabilities** – Basic and advance concepts of Statistics and Probability, Inferential Statistics - Python, R, SAS.
- ✓ **Analytics & Visualization** – Interactive visual Presentation of data and findings - Python, R, MATLAB, SAS, Tableau or other visualization tools.
- ✓ **Machine Learning** – To Apply the algorithm on data - Python, R, SAS & other IDEs like Microsoft Azure ML Studio, IBM SPSS, SAP Predictive Analytics.

++ Business or Domain Knowledge ++

What are the opportunities?

“Data Science is not a Technology but a whole New world”

- **Data Engineer**
- **Dashboard/Analytics Expert**
- **Data Analyst**
- **Machine Learning Expert**
- **Data Scientist – One for ALL!!**

Prerequisites & Target Audience

Prerequisite or who is suitable to get into Data Science?

Most Suitable : Professional already working as Data Analyst, Analytics Expert, BI Developers/Architects, ETL/Big data Engineer/Architects, DWH Designer/Architects.

Minimum Criteria :

- ✓ Basic understanding of programming concepts like PL/SQL, C, C++, Java & RDBMS.
- ✓ Basic knowledge of Mathematics and Statistics Concepts.
- ✓ Basic Knowledge of Reporting or Visualization like Tableau, Spotfire, SAP Business Objects or any other reporting tools.

Great Determinations! Dedications! Consistency! and Confidence !!

Get Set Ready!!!

Homework

- Do some more Research on What is Data Science & Machine Learning
- Try to understand CRISP-DM Methodology

<https://www.the-modeling-agency.com/crisp-dm.pdf>

“Qs & As”



THANK YOU

Data Science with Python

Day 2 - Data Science Project Life Cycle

Today's Agenda

✓ Data Science Project Life Cycle

- CRISP - DM Process Model
- CRISP - DM Phases
 - Business understanding
 - Data understanding
 - Data preparation
 - Modeling
 - Evaluation
 - Deployment

✓ Basics of Statistics

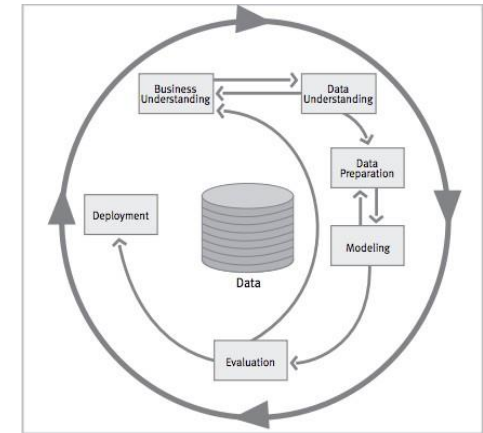
- Random Variable
- Type of Random Variables
- Central Tendencies
 - Mean
 - Mode
 - Median
- Probability, Probability Distribution of Random Variables

Data Science Project Life Cycle

What is CRISP-DM ?

CRISP-DM was conceived in late 1996 by three veterans of the young and immature data mining market. CRISP Stands for “**C**ross-**I**ndustry **S**tandard **P**rocess for **D**ata **M**ining”

This Process model for data mining provides an overview of the life cycle of a data mining project. It contains all the phases of a project, their respective tasks, and the relationships between these tasks. Relationships could exist between any data mining tasks depending on the goals, the background, and the interest of the user—and most importantly—on the data.



Data Science Project Life Cycle

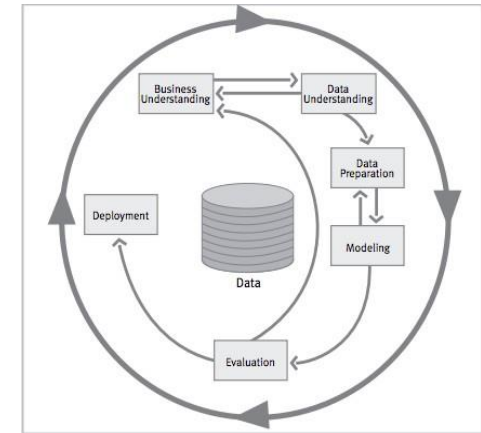
CRISP-DM Phases

The life cycle of a data mining project consists of six phases as shown –

1. Business understanding - This initial phase focuses on understanding the project objectives and requirements from a business perspective, then converting this knowledge into a data mining problem definition and a preliminary plan designed to achieve the objectives.

2. Data understanding - The data understanding phase starts with initial data collection and proceeds with activities that enable you to become familiar with the data.

3. Data preparation - The data preparation phase covers all activities needed to construct the final dataset or data that will be fed into the modeling tool(s) from the initial raw data. Tasks include table, record, and attribute selection, as well as transformation and cleaning of data for modeling tools.

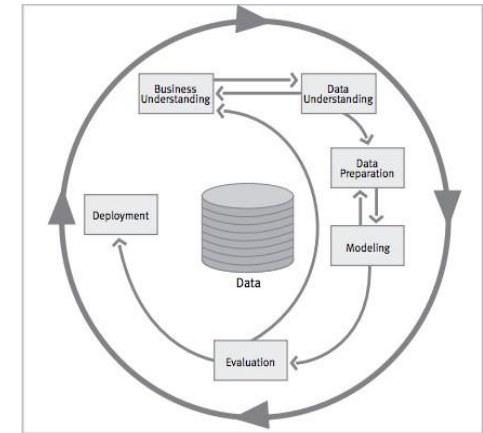


Data Science Project Life Cycle

CRISP-DM Phases

4.Modeling - In this phase, various modeling techniques are selected and applied, and their parameters are calibrated to optimal values.

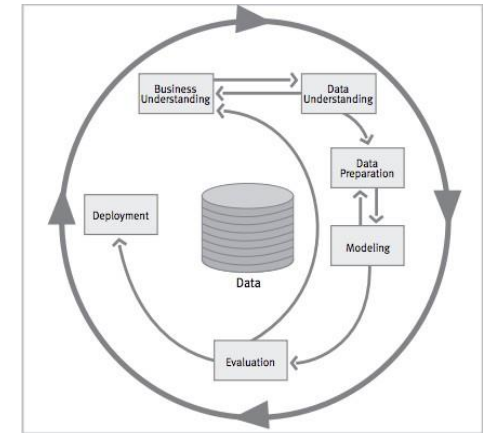
5.Evaluation - At this stage you have built a model (or models) that appears to have high quality from a data analysis perspective. Before proceeding to final deployment of the model, it is important to thoroughly evaluate it and review the steps executed to create it, to be certain the model properly achieves the business objectives. A key objective is to determine if there is some important business issue that has not been sufficiently considered.



Data Science Project Life Cycle

CRISP-DM Phases

6. Deployment- Creation of the model is generally not the end of the project. Even if the purpose of the model is to increase knowledge of the data, the knowledge gained will need to be organized and presented in a way that the customer can use it. It often involves applying “live” models within an organization's decision making processes. 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



Data Science Project Life Cycle

CRISP-DM Phases

Business Understanding	Data Understanding	Data Preparation	Modeling	Evaluation	Deployment
Determine Business Objectives <i>Background Business Objectives Business Success Criteria</i>	Collect Initial Data <i>Initial Data Collection Report</i>	Select Data <i>Rationale for Inclusion/ Exclusion</i>	Select Modeling Techniques <i>Modeling Technique Modeling Assumptions</i>	Evaluate Results <i>Assessment of Data Mining Results w.r.t. Business Success Criteria Approved Models</i>	Plan Deployment <i>Deployment Plan</i>
Assess Situation <i>Inventory of Resources Requirements, Assumptions, and Constraints Risks and Contingencies Terminology Costs and Benefits</i>	Describe Data <i>Data Description Report</i>	Clean Data <i>Data Cleaning Report</i>	Generate Test Design <i>Test Design</i>	Review Process <i>Review of Process</i>	Plan Monitoring and Maintenance <i>Monitoring and Maintenance Plan</i>
Determine Data Mining Goals <i>Data Mining Goals Data Mining Success Criteria</i>	Explore Data <i>Data Exploration Report</i>	Construct Data <i>Derived Attributes Generated Records</i>	Build Model <i>Parameter Settings Models Model Descriptions</i>	Determine Next Steps <i>List of Possible Actions Decision</i>	Produce Final Report <i>Final Report Final Presentation</i>
Produce Project Plan <i>Project Plan Initial Assessment of Tools and Techniques</i>	Integrate Data <i>Merged Data</i>	Format Data <i>Reformatted Data</i>	Assess Model <i>Model Assessment Revised Parameter Settings</i>		Review Project <i>Experience Documentation</i>
	Verify Data Quality <i>Data Quality Report</i>	Dataset <i>Dataset Description</i>			

Data Science with Python

Day 2 – Basics Of Statistics

Basics Of Statistics

What is Statistics?

Statistics is a branch of Mathematics dealing with the collection, analysis, interpretation, presentation, and organization of data.



Random variable

A Variable which is used to store value corresponding to each outcome of a Random Experiment/Event/Activity. Ex. Coin Flip, $RV = \{H, T\}$

Type of Random variable

Based on the nature of outcome RV can be Discrete or Continues.

Discrete - Finite measurement, no decimals, Ex. Number of people

Continues - Infinite measurements between 2 consecutive values, Ex. Weight, Age



Central Tendencies- Mean Mode, Median

Type	Description	Example	Result
Mean	Sum of values of a data set divided by number of values:	$(1+2+2+3+4+7+9) / 7$	4
Median	Middle value separating the greater and lesser halves of a data set	1, 2, 2, 3, 4, 7, 9	3
Mode	Most frequent value in a data set	1, 2, 2, 3, 4, 7, 9	2

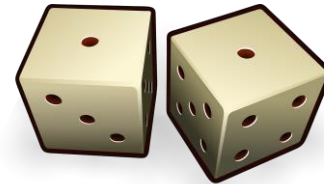


Basic Probability

Probability is the measure of the likelihood that an event will occur. In case of Random variable we are interested in knowing the Probabilities of getting different values.

Probability - $P(X=X_i) = F(X_i) / F_{\text{Total}}$

Ex. Rolling Dices, $P(\text{Output}=1) = 1/6$





Probability Distribution of RV

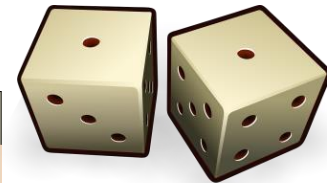
Table/ Chart/Formula to show relationship between Values and Corresponding Probabilities or shows the distribution of probabilities by values.

Type of Probability Distribution

Based on type of RV, Probability Distribution can also be either Discrete or Continuous

1. Discrete Probability Distribution – Probability Mass Function (PMF)

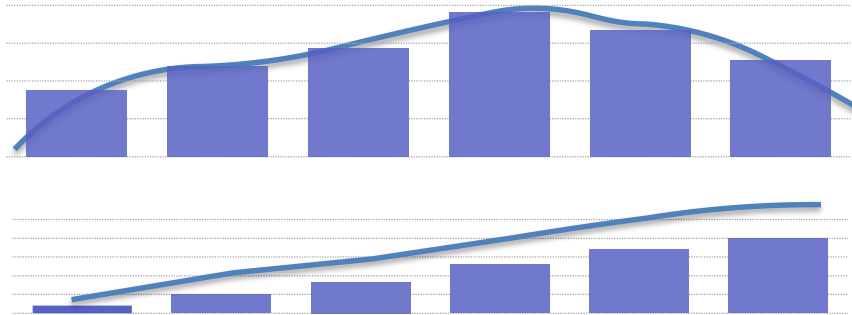
Discrete Probability Distribution (PMF)						
Values(X)	1	2	3	4	5	6
P(X)	1/6	1/6	1/6	1/6	1/6	1/6





Type of Probability Distribution

2. Continuous Probability Distribution – Probability Density Function (PDF)



AGE GROUP	#PERSONS(X) (IN K)	P(X)	CP(X)
0-10	11	0.105	0.105
10-20	15	0.143	0.248
20-30	18	0.171	0.419
30-40	24	0.229	0.648
40-50	21	0.200	0.848
50-60	16	0.152	1
Total	105	1	

“Qs & As”



THANK YOU