# Introduction to data science

LECTURE – 1 (WEEK 1)

# Overview

Course and Instructor's introduction

Understanding data

Data science intro
◦ The process

Examples

Types of Data

# Welcome to Intro to Data Science – fall 2022

Instructor (s):

Dr. Asma Ahmad [Sec: A, B, C]

◦ Email:  asma.ahmad@nu.edu.pk/asma.ahmad@lhr.nu.edu.pk

◦ Office: E-157

Office Hours: will be displayed on office door shortly

Course material on Google classroom:  <span style="color:red">xv2e6ny</span>

TA:

◦ TBD

# Reference Books and Resources

**Cathy O'Neil and Rachel Schutt. Doing Data Science, Straight Talk From The Frontline. O'Reilly. 2014. ISBN 978-1-449-35865-5.**

**Jiawei Han, Micheline Kamber and Jian Pei. Data Mining: Concepts and Techniques, Third Edition. Morgan Kaufmann Publishers. 2012. ISBN 978-0-12-381479-1.**

**Tom Mitchell: Machine Learning**

Jure Leskovek, Anand Rajaraman and Jeffrey Ullman. Mining of Massive Datasets. v2.1, Cambridge University Press. 2014.

Kevin P. Murphy. Machine Learning: A Probabilistic Perspective. MIT Press. 2013. ISBN 0262018020.

Foster Provost and Tom Fawcett. Data Science for Business: What You Need to Know about Data Mining and Data-analytic Thinking. O'Reilly 2013. ISBN 978-1-449-36132-7.

# Tools and Software Packages

WEKA

KNIME

ORANGE

CBA

MATLAB

PYTHON

R

SPSS

SAS

ArcView GIS

Maptitude for the Web

etc.,

**Language:  Python**
**Tools: NumPy, Pandas, matplotlib, scikit-learn**
**Other packages if needed: SciPy, and SymPy**

# Tentative Grading Policy

| Class | |
|---|---|
| Assignments | 15% |
| Quizzes | 10% |
| Midterm Exam | 30% |
| Final Exam | 45% |

# Grading Policy

There is simply no chance of extension in any of the deadline, what so ever

You can request for re-checking of any of your evaluation as per following rules;

Exams: <span style="color:green">Same day</span>

Assignments: <span style="color:green">2 days after handing-over</span>

Quizzes: <span style="color:green">2 days after handing over</span>

<span style="color:red">Warning! After due time, request will not be considered even if it's genuine.</span>

<span style="color:red">Warning! Regularly check flex and don't come in the end with bulk of queries in hand.</span>

# General Guidelines

Visit Google classroom regularly for updates

No email submissions when Google classroom is there. *Always remember that, you are putting your task in trash by yourself when you are emailing it.*

❑ Cheating cases are intolerable. You will be given negative marks for cheated stuff irrespective of the fact that, you were provider or the other one. Your cheating in exam will make it easy for you to step down from the Course with an 'F' grade. . . ☹

❑ There will be no re-take of any evaluation if you haven't informed earlier through a proper channel.

Quiz is inevitable so always, expect a One ☺ [at least, one in every week]

# General Guidelines

You have to depend on yourself to have a good grade in the course

◦ *You need to appear in demo to have it graded otherwise you will be given 5% of the total marks even if it was best assignment/project of the class*

◦ *Grading will be individual even for group tasks.*

◦ *There will always be a quiz of written assignments whose performance will be considered as performance in that assignment.*

# Tentative Outline

**Introduction & Applications**

**Data (Acquisition, Storage, & processing)**

**Data Preprocessing and Mining**
- **Classification**
- **Clustering**
- **Association Rule Mining**
- **Attribute selection (Feature Selection)**

**Machine Learning**
- **Classification and Regression**
- **Gradient Descent**
- **Regularization**
- **Support Vector Machines**

- **Dimensionality Reduction**
- **Outlier Detection**

**Statistical Inference**
- **Statistical Modeling**
- **Probability distribution**
- **Fitting a model**

**Descriptive and Exploratory Analysis**

**Market Basket Analysis: Market basket analysis is a data mining technique used by retailers to increase sales by better understanding customer purchasing patterns**

**Exploratory Data Analysis [EDA]**

**Data Visualization**

# Learning Outcomes

Understand the basics of Data Science,

Prepare and wrangle the data for analysis

Perform Exploratory Data Analysis (EDA)

Understand and apply machine learning algorithms to gain insight from the data

# Why do we need Data Science?

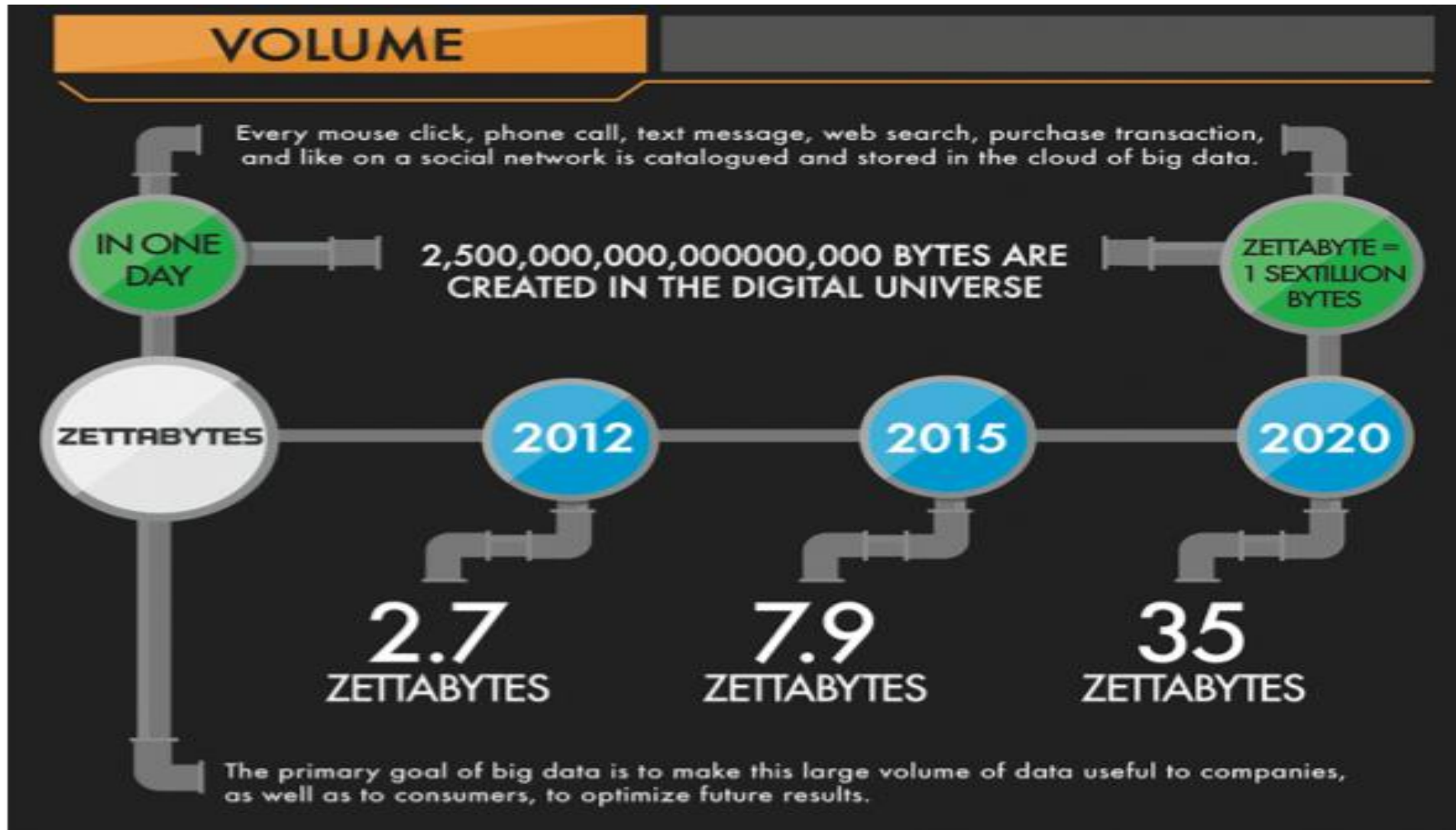Good decisions require good information derived from raw facts

Big Data

"Between the dawn of civilization and 2003, we only created five exabytes of information; now we're creating that amount every two days."

- Eric Schmidt, Google

# Why do we need Data Science?

# Why do we need Data Science

# Nate Silver

American Data Scientist who analyzes elections and baseball.
 -PECOTA: a system for forecasting the performance and career development of Major -League Baseball players.
- 2012 U.S. Presidential election

# 2012 U.S. Presidential election

o Correctly predicted the winners of all the states.
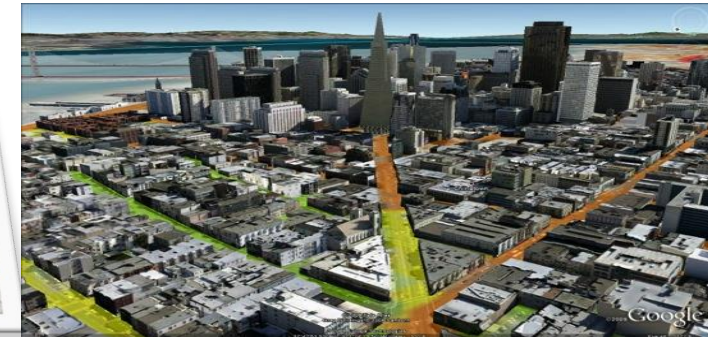


"Nate Silver won the election"
– Harvard Business Review

538 prediction | actual

# What can you do with the data?



Crowdsourcing + physical modeling + sensing + data assimilation

to produce:

From Alex Bayen, UCB

# Data vs. Information

| Data | Information |
|------|-------------|

**Raw facts**
- Have not yet been processed to reveal their meaning to the end user

**Building blocks of information**

**Data management**
- Generation, storage, and retrieval of data

Produced by processing raw data to reveal its meaning

Requires context

Bedrock of **knowledge**

Should be accurate, relevant, and **timely** to enable good decision making

# Data vs. Information (cont'd.)

**FIGURE 1.1** Transforming raw data into information



### a) Data entry screen

### b) Raw data

### c) Information in summary format

| Rank | COUNT | %/INFS | TOT/COL | %/COL. TOT. | %/COL. FAC. |
|---|---|---|---|---|---|
| Adjunct | 5 | 20.00% | 23 | 21.74% | 3.27% |
| Assistant Professor | 2 | 8.00% | 28 | 7.14% | 1.31% |
| Associate Professor | 9 | 36.00% | 37 | 24.32% | 5.88% |
| Instructor | 2 | 8.00% | 18 | 11.11% | 1.31% |
| Professor | 7 | 28.00% | 47 | 14.89% | 4.58% |

### d) Information in graphical format

SOURCE: Course Technology/Cengage Learning
Data entry screen courtesy of Sedona Systems, 2011.
Information screens courtesy of JCBDashboard, 2011.

# Data, Information, and Beyond

# What is Data Science?

Data Science is the application of computational and statistical techniques to address or gain insight into some problem in the real world

Data Science =  Statistics +  data processing + machine learning + scientific inquiry + visualization + business analytics + big data + …

# Contrast: Databases

| Datawarehouse | Data Science |
|---|---|
| Querying the past | Querying the future |



**Business intelligence** (**BI**) is the transformation of raw data into meaningful and useful information for business analysis purposes. BI can handle enormous amounts of unstructured data to help identify, develop and otherwise create new strategic business opportunities - Wikipedia

INTRO TO DS

# Big Data – Five V of Data

Volume:
- ◦ How much data is really relevant to the problem solution? Cost of processing?
- ◦ *So, can you really afford to store and process all that data?*

Velocity:
- ◦ Much data coming in at high speed
- ◦ Need for streaming versus block approach to data analysis
- ◦ *So, how to analyze data in-flight and combine with data at-rest*

Variety:
- ◦ A small fraction is structured formats, Relational, XML, etc.
- ◦ A fair amount is semi-structured, as web logs, etc.
- ◦ The rest of the data is unstructured text, photographs, etc.
- ◦ *So, no single data model can currently handle the diversity*

# Big Data – Five V of Data (Cont.)

Veracity: cover term for …
- Accuracy, Precision, Reliability, Integrity
- *So, what is it that you don't know you don't know about the data?*

Value:
- How much value is created for each unit of data (whatever it is)?
- *So, what is the contribution of subsets of the data to the problem solution?*

# Types of Data Analytics

***Descriptive***: A set of techniques for reviewing and examining the data set(s) to understand the data and analyze business performance.

***Diagnostic***: A set of techniques to determine what has happened and why

***Predictive***: A set of techniques that analyse current and historical data to determine what is most likely to (not) happen

# Types of Data Analytics (Cont.)

***Prescriptive***: A set of techniques for computationally developing and analyzing alternatives that can become courses of action – either tactical or strategic – that may discover the unexpected

***Decisive***: A set of techniques for visualizing information and recommending courses of action to facilitate human decision-making when presented with a set of alternatives.

# Data Science Life Cycle

**Problem Understanding**: It all starts with understanding the problem at hand, the questions, and the answers we are trying to find from the dataset at hand.

**Data Acquisition**: Data Acquisition, as the name suggests, is about retrieving the data with the help of Data Engineers where required. It also consolidates all of the data required to answer the question or to solve the problem at hand.

**Data Wrangling**: Data wrangling is about using knowledge to preprocess data. It involves looking for missing values and asking business questions like why they are missing. Furthermore, it uses knowledge to give shape to the dataset appropriate for visualizations and to support the coming steps in the life cycle.

**Data Exploration**: Data Exploration is about visualization and other statistics' measures to see whether the questions we asked, in the beginning, are being answered or not? The data analyst's job ends here.
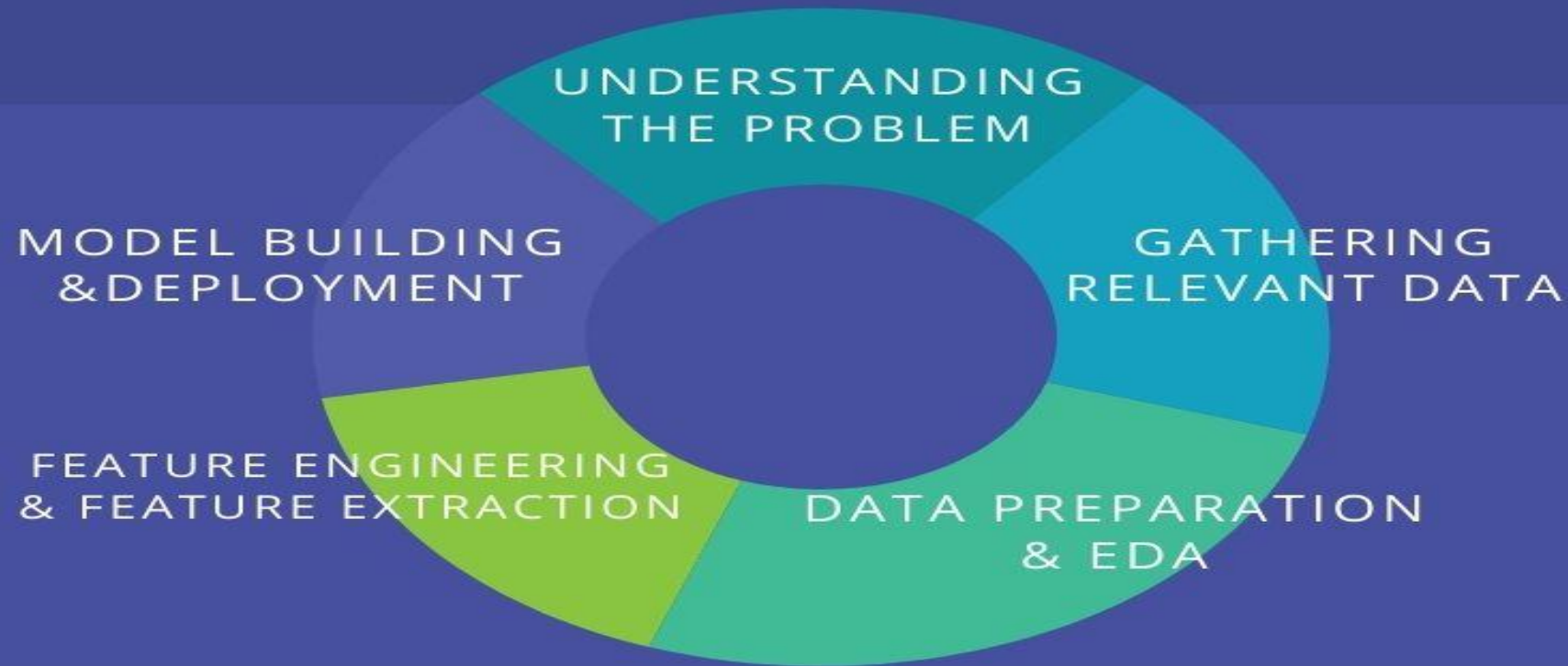
# Data Science Life Cycle

**Feature Engineering and Selection**: It is a preprocessing step before modeling in both Machine Learning and Deep Learning. We will look into these fields in the coming sections. It has similar steps to Data Wrangling apart from some algorithms for Feature Selection and transformation.

**Modeling**: Modeling is the process that uncovers the meaning of the data. It is about capturing underlying trends and the data's behavior to make the model, which can be used for predictive analytics as described in the previous section.

**Deployment**: After we build the model we'll deploy it in the most efficient and optimized manner so that real-world people can use it. It can be deployed on mobile applications and web applications.

**Monitoring**: After we have deployed the model, we will want to monitor it. Monitoring is about familiarizing the model with the new dataset and tracking the number of requests that the model receives. It also involves making changes to the analysis and starting over if required.

# Asking Good Questions

Software developers are not encouraged to ask questions, but data scientists are:
◦ What exciting things might you be able to learn from a given data set?
◦ What things do you/your people really want to know?
◦ What data sets might get you there?

e.g., Baseball
◦ How to best measure individual player's skill, value or performance?
◦ How fair do trades between teams work out?
◦ What is the trajectory of player's performances as they mature and age?
◦ To what extent does batting performance correlate with the position played?

# Structured vs. Semi-Structured vs. Unstructured Data

**Structured Data**

It comes with a predefined format and structure. Structured Data is usually stored in Relational Databases. It is easy to deal with in the Data Science domain.

| Sepal_length | Sepal_width | Petal_length | Petal_width | Species |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | versicolor |
| 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 4.6 | 3.1 | 1.5 | 0.2 | virginica |

# Semi-Structured Data

It comes with a predefined format and structure but is not stored in the Relational Database.

- JSON (Javascript Object Notation)

```json
{
    name: "Linear Algebra for Machine Learning",
    author: "Json Brownlee",
    pages: 211,
    parts: 5,
    format: "PDF",
    total_codes: 92
}
```

- XML (Extensible Markup Language)

```xml
<?xml version="1.0" encoding="UTF-8"?>
<book>
    <author>Json Brownlee</author>
    <format>PDF</format>
    <name>Linear Algebra for Machine Learning</name>
    <pages>211</pages>
    <parts>5</parts>
    <total_codes>92</total_codes>
</book>
```

# Unstructured Data

It does not have a specific format and lacks structure. It is the type of data that presents many challenges to handle in the Data Science domain

**Examples**:

◦ Images

◦ Videos

◦ Speech

# THAT'S IT FOR TODAY!