# Foundations of AI & Data Science

Unit 1

# LEARNING OBJECTIVES

- To excite students about the potential that resides in data and the value that data analytics can add to business processes

- To impart skills related to data cleaning/wrangling, data transformation/preprocessing, and data comprehension through statistical analysis

- To impart skills related to analytical (mathematical) data modeling

# Course overview

- Introduction to Data Science & AI

- Data Preparation and Cleaning

- Classification Techniques

- Supervised & Unsupervised Learning Methods

- Data and Customer Segmentation

- Time Series Forecasting

# Software and data repositories

- Python

- Data on Kaggle Website
  - http://www.kaggle.com/

# Books

- Data Science for Business by Provost and Fawcett (2013)

- Data Mining Concepts and Techniques by Han and Kamber (2011)

# MOOCs on Data Sciences

- In the past couple of years, Data Science related courses and specialization have been extremely popular on MOOCs websites:

  - John Hopkins University (Coursera)

  - University of Washington (Coursera)

  - Google (UdaCity)

  - UC Berkley (EdX)

  - University of Toronto (Coursera)

  - And many others……..

# Artificial Intelligence

Origin and
Taxonomy

# Data is Everywhere!

- Lots of data is being collected and warehoused
  - Web data
    - Google has Peta Bytes of web data
    - Facebook has billions of active users
  - purchases at department/ grocery stores, e-commerce
    - Amazon handles millions of visits/day
  - Bank/Credit Card transactions
- Competitive Pressure is Strong
  - Provide better, customized services

# What Is Big Data?

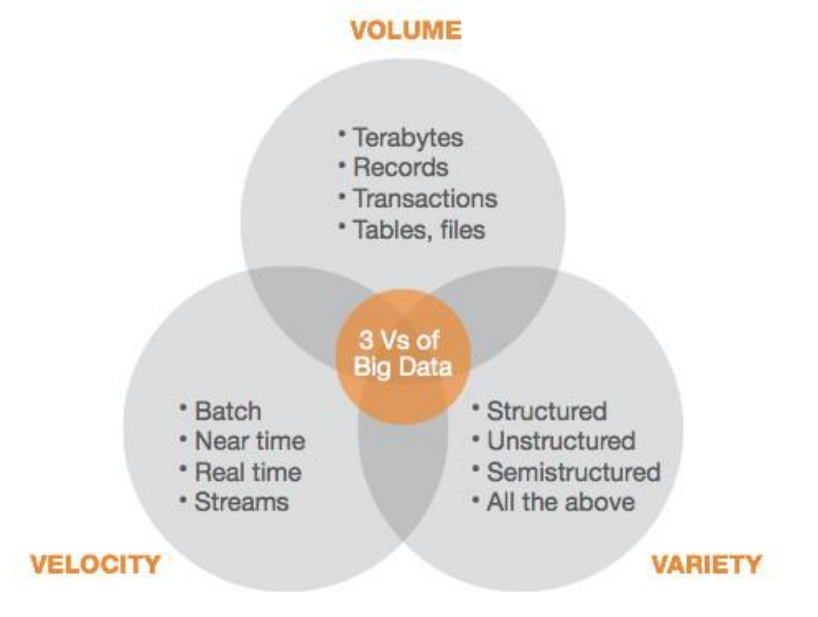- There is not a consensus as to how to define big data

"Big data refers to data sets whose size is beyond the ability of typical database software tools to capture, store, manage and analyze."
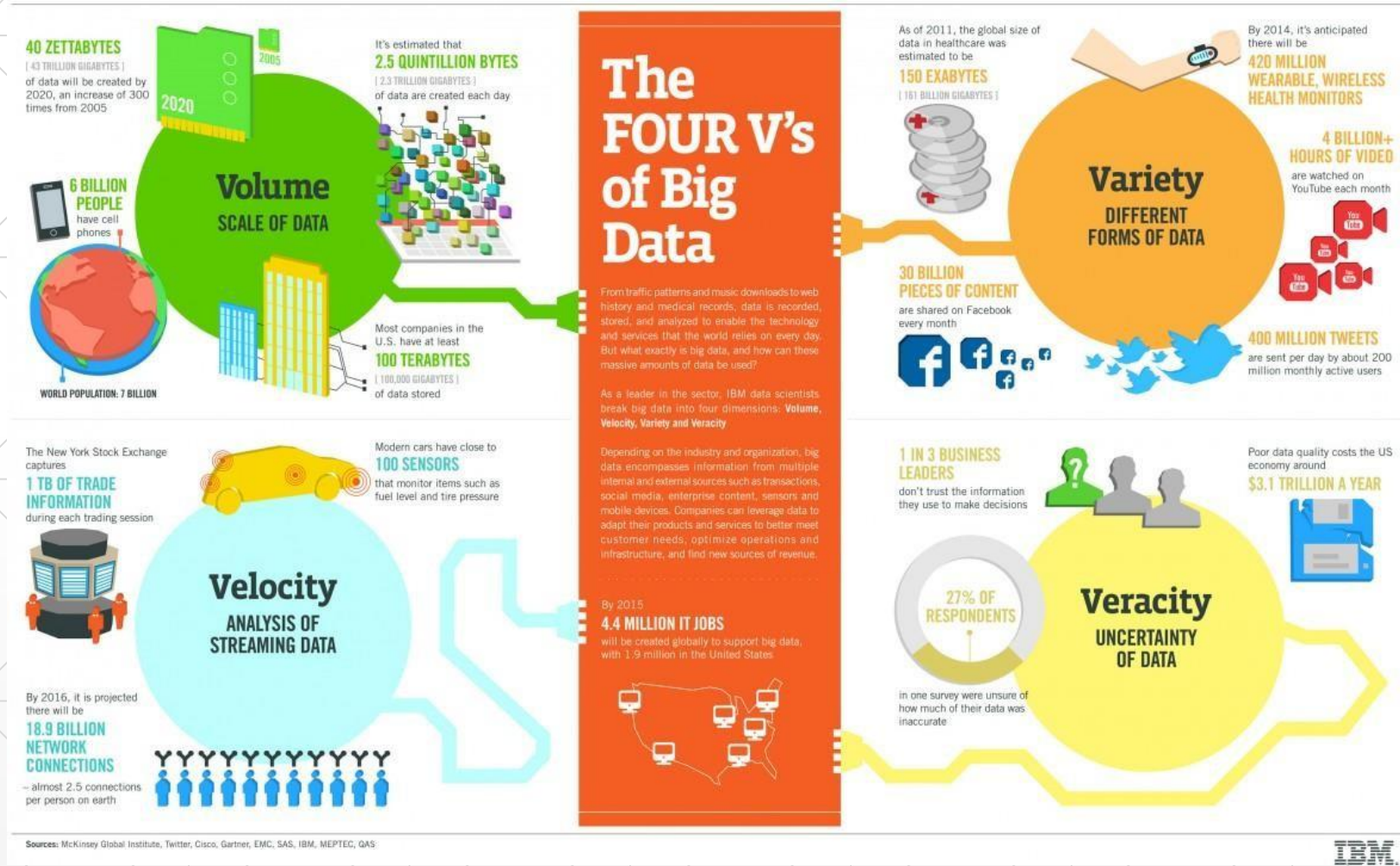
*- The McKinsey Global Institute, 2011*

- One reasonable definition is that it's data which can't comfortably be processed on a single machine.

# 3 V's of Big Data

- Doug Laney was the first one in talking about 3 V's in Big Data management:

  - **Volume**: there is more data than ever before, its size continues increasing, but not the percent of data that our tools can process

  - **Variety**: there are many different types of data, as text, sensor data, audio, video, graph, and more

  - **Velocity**: data is arriving continuously as streams of data, and we are interested in obtaining useful information from it in real time



VOLUME
- Terabytes
- Records
- Transactions
- Tables, files

3 Vs of Big Data

VELOCITY
- Batch
- Near time
- Real time
- Streams

VARIETY
- Structured
- Unstructured
- Semistructured
- All the above

# 4 V's by IBM (2014)

# What to do with Data?

**Health Care**
- Disease Prediction
- Drug Discovery
- Virtual Assistant
- Image Analysis

**Banking**
- Fraud Detection
- Credit Risk Modeling

**Ecommerce**
- Produc Recommendation
- Fake review detection

**Transport**
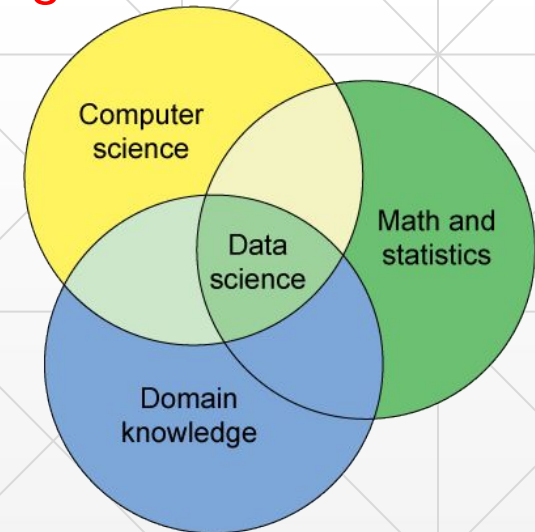- Self Driving Cars
- Traffic control

# Data Science

- Data science incorporates varying elements and builds on techniques and theories from many fields, including ***mathematics, statistics, data engineering, pattern recognition and learning, advanced computing, visualization, uncertainty modeling, data warehousing, and high performance computing*** with the goal of extracting meaning from data and creating data products. *(From Wikipedia)*

- A practitioner of data science is called a data scientist.

# Why Data Science?

- Traditional Data – Analyzed using simple BI tools
- Unstructured data is generated from different sources like <span style="color:red">financial logs, text files, multimedia forms, sensors, and instruments</span>.
- Huge volume and variety of data requires - Complex and advanced analytical tools and algorithms for meaningful insights out of it.



**Unstructured data** will account for **more than 80%** of the data collected by organizations

UNSTRUCTURED DATA

STRUCTURED DATA

Total Data Stored

1980  1990  2000  2010  2020

Source: https://www.edureka.co/blog/what-is-data-science/