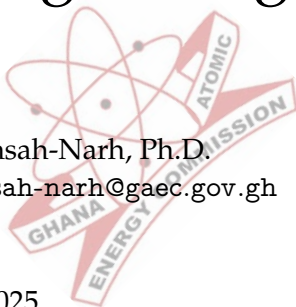Ghana Space Science & Technology Institute, GAEC

# Classification in Deep Learning: A Beginner's Guide

Theophilus Ansah-Narh, Ph.D.
theophilus.ansah-narh@gaec.gov.gh

February 19, 2025

1. Big Data

2. Machine Learning

3. Deep Learning

4. Performance Metrics

5. DEMO

6. Remarks

- **Importance of Bankruptcy Prediction:**
  - Accurately evaluating a company's financial health is crucial for stakeholders to mitigate risks and prevent bankruptcy.
  - Effective bankruptcy prediction helps stakeholders make informed decisions and proactively manage financial stability in a rapidly changing business landscape.
- **Overview of the Study's Methodology:**
  - The study presents an innovative approach by integrating Domain Adaptation Learning (DAL) and Genetic Algorithm (GA) techniques.
  - DAL addresses distributional changes in real-world scenarios, while GA excels in feature selection.
  - Six machine learning models are rigorously evaluated against the hybrid model to enhance corporate bankruptcy prediction.

# Big Data

**Facets and Elements of Big Data.** *Image credit*: Dzone website

- **Big data** is an amount of data that is enormous in volume and is constantly expanding rapidly.
- No typical data management systems can effectively store or analyze this data because of its magnitude and complexity.
- fundamental characteristics of big data are listed below

Volume

# The Erra of Big Data III

- Big Data is a vast *volume* of data generated from many sources daily, such as business processes, machines, social media platforms, networks, human interactions, and many more.

- Industry trends predict a significant increase in data volume over the next few years.

- Usually measured in gigabytes (GB), terabytes (TB), zettabytes (ZB), and yottabytes (YB)

- Nonetheless, Big data generally refers to datasets with a high volume of the order of magnitude of exabytes ($10^{18}$B $= 10^9$GB $= 10^6$TB $= 1$EB ) and greater (Jelic *et* al. 2019).

**DAILY ACTIVE USERS** 1.96 Billion

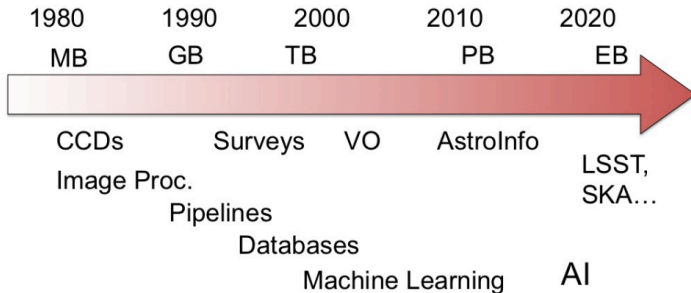**DATA GENERATED EVERY DAY** 500+ TERABYTES

10M+ GROUPS

300M+ STORY UPDATES

*Image credit*: Analytics Vidhya

# The Evolving Data-Rich Astronomy

An example of a "Big Data" science driven by the advances in computing/information technology



**Key challenges: data heterogeneity and complexity**
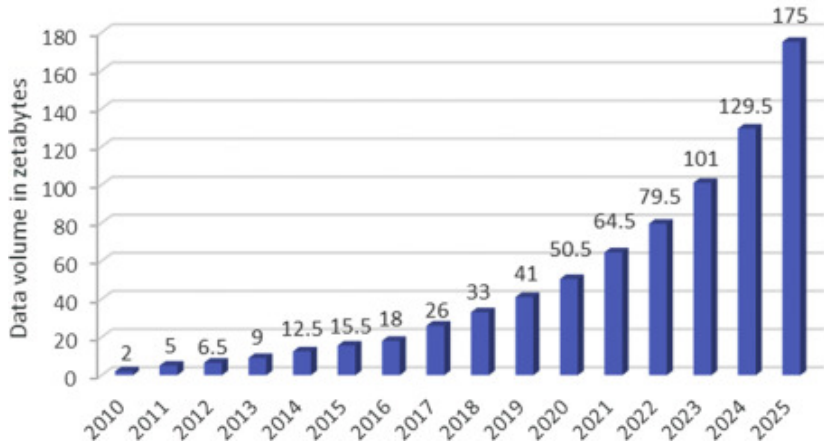
*Image credit*: Djorgovski, 2019

## Global data volume



*Image credit*: Lei & Kong, 2020

## Variety

- In the past, data is only collected from **databases** and **sheets**.
- These days the data will come in array forms, that are **PDFs**, **Emails**, **audios**, **SM posts**, **photos**, **videos**, etc.
- Big Data can be **structured**, **unstructured**, and **semi-structured** that are being collected from different sources.

- **Structured data:** In Structured schema, along with all the required columns. It is in a tabular form. Structured Data is stored in the relational database management system.
- **Semi-structured:** In Semi-structured, the schema is not appropriately defined, e.g., JSON, XML, CSV, TSV, and email. OLTP (Online Transaction Processing) systems are built to work with semi-structured data. It is stored in relations, i.e., tables.
- **Unstructured Data:** All the unstructured files, log files, audio files, video files, e-mails, word processing, and image files are included in the unstructured data.

## Veracity

- The accuracy of your findings can be severely harmed by poor data reliability.
- Making it one of the most crucial big data qualities
- There's a need to calibrate your data since most of the data you encounter is unstructured.
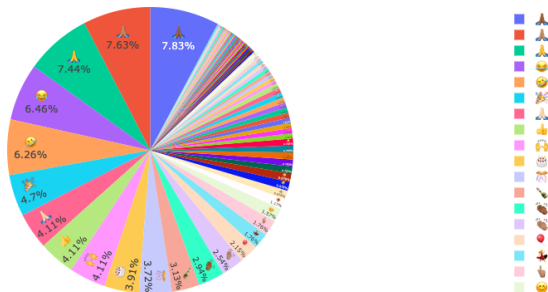
```
In [13]: print('Group wise Stats')
         print("Messages:", total_messages)
         print('Media:', media_messages)
         print('Emojis:', emojis)
         print('Links:', links)
```
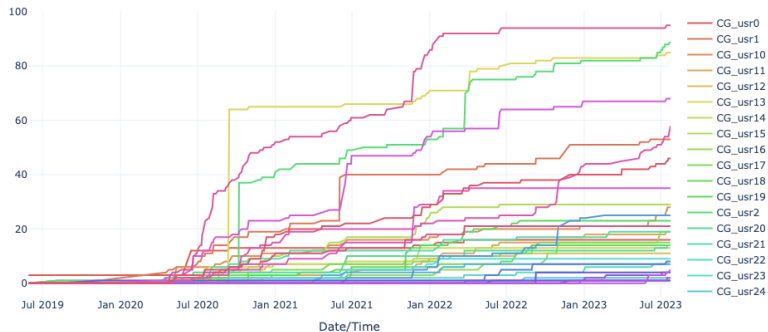
```
Group wise Stats
Messages: 845
Media: 182
Emojis: 511
Links: 188
```
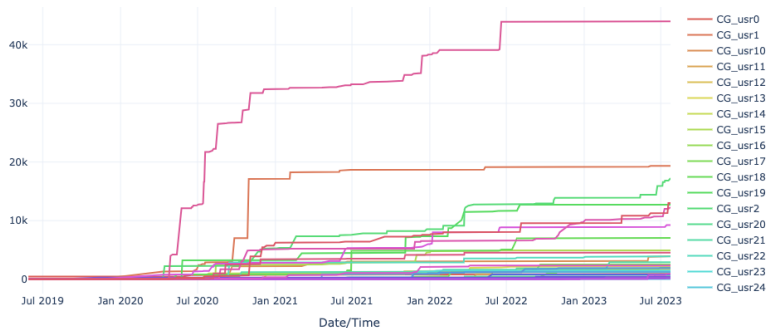
Emoji Distribution

User interventions count (cumulative)

Count of sent characters (cumulative)

Response matrix
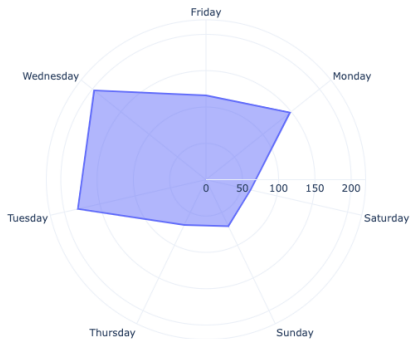
## Value

- On this data set, analysis and pattern recognition are performed.
- The results of the method may be used to determine the value of the data.
- Making it one of the most crucial big data qualities.

# Machine Learning

*Image credit*: mapendo site

*Image credit*: Mehra & Hasanuzzaman, (2020)

T. Ansah-Narh

TOOLBOX GSSTI-GAEC

# Deep Learning

*Image credit*: Odi & Nguyen, (2018)

Schematic of a feed-forward neural network

*Image credit*: Zhu *et* al , (2019)

Image

Convolved
Feature

*Image credit*: Medium

Input

T. Ansah-Narh

GSSTI-GAEC

Output

# Performance Metrics

|            | Predicted 0 | Predicted 1 |
| ---------- | ----------- | ----------- |
| Actual 0   | TN          | FP          |
| Actual 1   | FN          | TP          |

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | TN | FP |
| Actual 1 | FN | TP |

| Metric | Formula | Evaluation focus |
|---|---|---|
| Accuracy | $\text{ACC} = \dfrac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$ | Overall effectiveness of a classifier |
| Precision | $\text{PRC} = \dfrac{\text{TP}}{\text{TP} + \text{FP}}$ | Class agreement of the data labels with the positive labels given by the classifier |
| Sensitivity | $\text{SNS} = \dfrac{\text{TP}}{\text{TP} + \text{FN}}$ | Effectiveness of a classifier to identify positive labels. Also called true positive rate (TPR) |
| Specificity | $\text{SPC} = \dfrac{\text{TN}}{\text{TN} + \text{FP}}$ | How effectively a classifier identifies negative labels. Also called true negative rate (TNR) |
| $F_1$ score | $F_1 = 2\,\dfrac{\text{PRC} \cdot \text{SNS}}{\text{PRC} + \text{SNS}}$ | Combination of precision (PRC) and sensitivity (SNS) in a single metric |
| Geometric mean | $\text{GM} = \sqrt{\text{SNS} \cdot \text{SPC}}$ | Combination of sensitivity (SNS) and specificity (SPC) in a single metric |
| Area under (ROC) curve | $\text{AUC} = \int_{0}^{1} \text{SNS} \cdot d\text{SPC}$ | Combined metric based on the receiver operating characteristic (ROC) space (*Powers, 2011*) |

# DEMO

# Remarks

*Image credit*: spiceworks site