# CIS 635 – 04 Knowledge Discovery & Data Mining
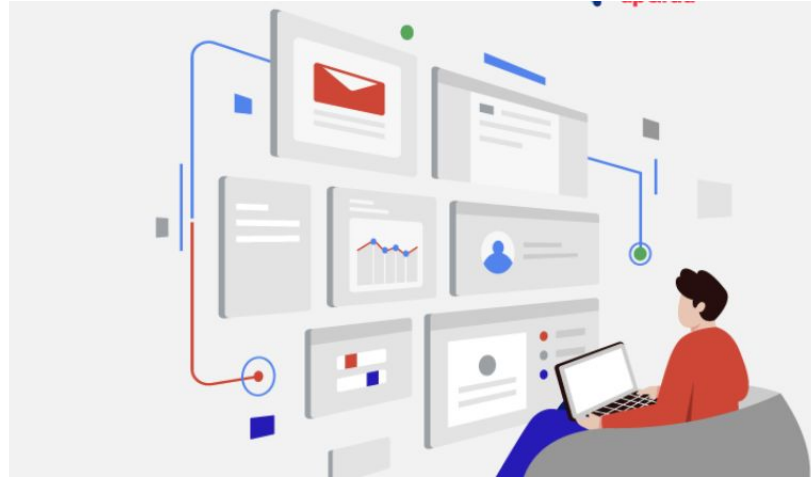
Data modalities - Introduction and Overview

# Outline

- Everyday observed data and their modalities and types
- Data Trends, and where DS fits

# Everyday observed (real life) data

# Everyday observed (real life) data



- We have seen an **explosion of data** in recent days, especially in last few years!
- The transition from the Analog to the Digital World has made this Big Shift (was unimaginable a few decades ago)

# Everyday observed (real life) data



- We have seen an **explosion of data** in recent days, especially in last few years!
- The transition from the Analog to the Digital World has made this Big Shift (was unimaginable a few decades ago)
- It's very difficult to **segregate them into disjoint types** as data come in different **format** and in different **modalities**.

# Everyday observed (real life) data



- We have seen an **explosion of data** in recent days, especially in last few years!
- The transition from the Analog to the Digital World has made this Big Shift (was unimaginable a few decades ago)
- It's very difficult to **segregate them into different types** as data come in different **format** and in different **modalities**.
- We will try a **simplified approach** to **group** them and try to cover most (NOT ALL) data we perceive in our everyday life.

# Everyday observed (real life) data

- Structured data
- Un(and Semi) structured data

# Structured data



Kaggle loan approval dataset

- Generally organized in **tables** and collected through filling **forms** (manual or online)
- Stored in **databases/spread-sheets** mainly
- Also popular the **.csv** file format

- Opening a bank account
- University registration
- Gmail
- Amazon account
- Your health profile

# Structured data



[Kaggle loan approval dataset](#)

- Generally collected through **forms** (manual or online)
- Stored in **databases/spread-sheets** mainly
- Also popular the **.csv** file format

- Opening a **bank account**
- **University registration**
- **Gmail, Azure, and/or Amazon** account
- Your **i**mmigration, health, social media **profile**

# Structured data



- There also other formats, or you can convert them to
- Some are stand alone, while others are sequences and/or series
- Software generated logs

- Genomics data
- Stock prices
- Your CC history
- Weather data
- Google maps

# Structured data



- There also other formats, or you can convert them to
- Software generated logs
- Some are stand alone, white others are sequences and/or series

- Genomics data
- Stock prices
- Your CC history
- Weather data
- Google maps

# Un/Semi Structured data









- Free forms
- Stored in data lake/warehouses mainly
- **Languages**: sequence of strings (semantics)
- **Audio**: Language + Acoustics; Music
- **Image** : Visual representation of the world
- **Video**: Sequence of images

- Some are sequence, while others are stand alone
- Social media data (emotions, vives)

# Un/Semi Structured data





- Free forms
- Stored in data lake mainly?
- Languages: sequence of strings
- Audio: Language + Acoustics; Music
- Image : Visual representation of the world
- Video: Sequence of images





- Some are sequences, while others are stand alone
- Social media data (emotions, vives)

# Un/Semi Structured data



- Free forms
- Stored in data lake mainly?
- Languages: sequence of strings
- Audio: Language + Acoustics; Music
- Image : Visual representation of the world
- Video: Sequence of images

- Some are sequence, while others are stand alone
- **Social media** data (discussions, messages, emotions, vives)

# Data Growth and the Trend

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

**Zetabyte** (y-axis)

| Year | Value |
|------|-------|
| 2010 | 2 |
| 2011 | 5 |
| 2012 | 6.5 |
| 2013 | 9 |
| 2014 | 12.5 |
| 2015 | 15.5 |
| 2016 | 18 |
| 2017 | 26 |
| 2018* | 33 |
| 2019* | 41 |
| 2020* | 64.2 |
| 2021* | 79 |
| 2022* | 97 |
| 2023* | 120 |
| 2024* | 147 |
| 2025* | 181 |

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- **Exponential growth**
  - *Before 2010 collectively relatively close to Zero*
- Mainly due to
  - IoT sensors
  - Social media
  -
- Forecast might be underestimated as we expect more due to the recent Generative AI/LLM releases (ChatGPT for an example)

## statista

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

**Zetabyte**

200

50

00

50

0

2 — 2010
5 — 2011
6.5 — 2012
9 — 2013
12.5 — 2014
15.5 — 2015
18 — 2016
26 — 2017
33 — 2018*
41 — 2019*
64.2 — 2020*
79 — 2021*
97 — 2022*
120 — 2023*
147 — 2024*
181 — 2025*

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- **Exponential growth**
  - *Before 2010 collectively relatively close to Zero*
- Mainly due to
  - IoT sensors
  - Social media
  - 
- Forecast might be underestimated as we expect more due to the recent Generative AI/LLM releases (ChatGPT for an example)

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

Zetabyte (chart)

2, 5, 6.5, 9, 12.5, 15.5, 18, 26, 33, 41, 64.2, 79, 97, 120, 147, 181

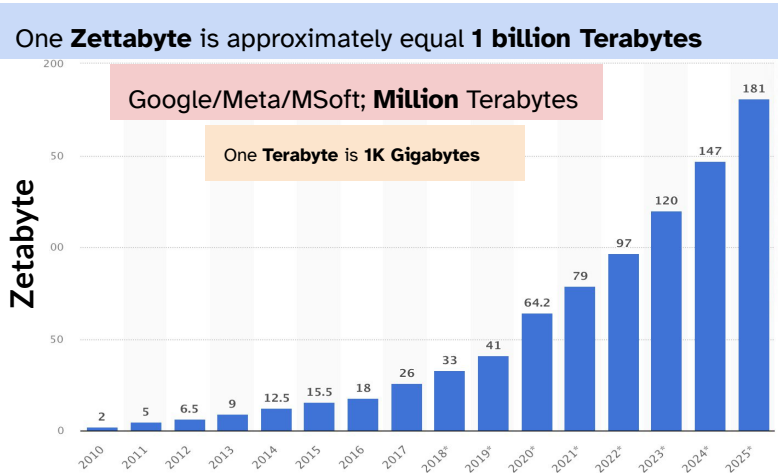2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018*, 2019*, 2020*, 2021*, 2022*, 2023*, 2024*, 2025*

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- Exponential growth
  - Before 2010 collected close to Zero
- **Mainly due to**
  - **Social media**
  - IoT sensors
  -
- Forecast might be underestimated as we expect more due to the recent Generative AI/LLM releases (ChatGPT for an example)

statista

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

**Zettabyte** (y-axis)

Chart values by year:
- 2010: 2
- 2011: 5
- 2012: 6.5
- 2013: 9
- 2014: 12.5
- 2015: 15.5
- 2016: 18
- 2017: 26
- 2018*: 33
- 2019*: 41
- 2020*: 64.2
- 2021*: 79
- 2022*: 97
- 2023*: 120
- 2024*: 147
- 2025*: 181

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- Exponential growth
  - Before 2010 coll... close to Zero
- **Mainly due to**
  - Social media
  - **IoT sensors**
  - 
- Forecast might be underestimated as we expect more due to the recent Generative AI/LLM releases (ChatGPT for an example)

THE INTERNET of THINGS
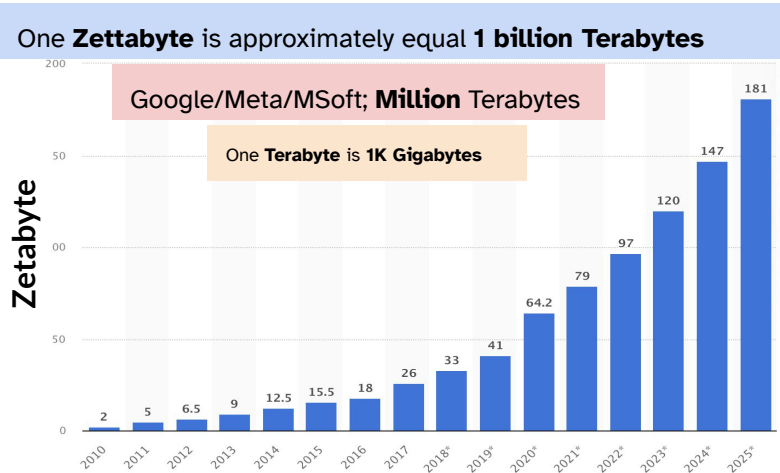
statista

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**



Zetabyte

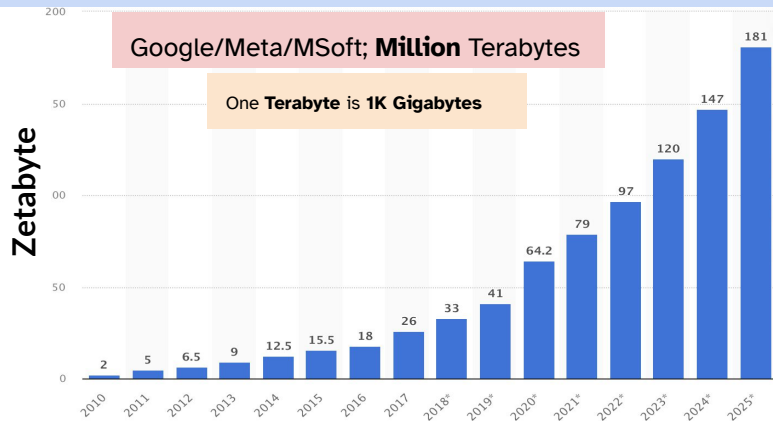| Year | Value |
|---|---|
| 2010 | 2 |
| 2011 | 5 |
| 2012 | 6.5 |
| 2013 | 9 |
| 2014 | 12.5 |
| 2015 | 15.5 |
| 2016 | 18 |
| 2017 | 26 |
| 2018* | 33 |
| 2019* | 41 |
| 2020* | 64.2 |
| 2021* | 79 |
| 2022* | 97 |
| 2023* | 120 |
| 2024* | 147 |
| 2025* | 181 |

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- Exponential growth
  - Before 2010  collectively relatively close to Zero
- Mainly due to
  - Social media
  - IoT sensors
  -
- Forecast might be underestimated as we expect more due to the recent Generative AI/LLM releases (ChatGPT for an example)
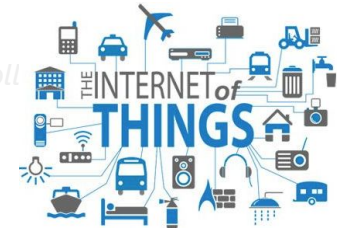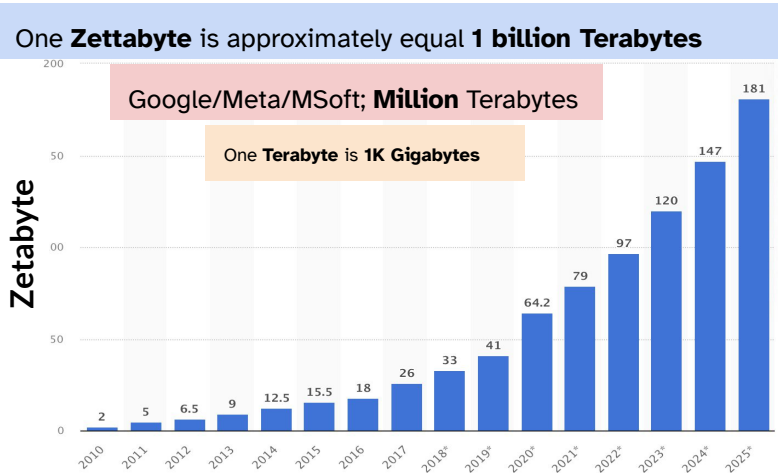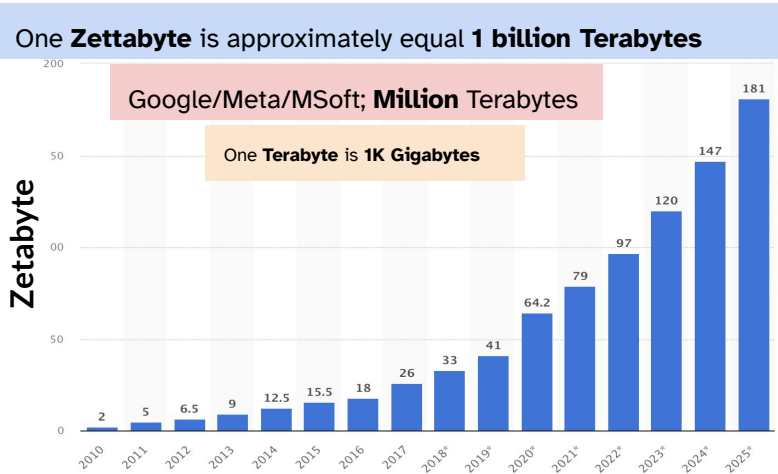
statista

# Data Growth and the Trend



One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

Zetabyte

| Year | Value |
|------|-------|
| 2010 | 2 |
| 2011 | 5 |
| 2012 | 6.5 |
| 2013 | 9 |
| 2014 | 12.5 |
| 2015 | 15.5 |
| 2016 | 18 |
| 2017 | 26 |
| 2018* | 33 |
| 2019* | 41 |
| 2020* | 64.2 |
| 2021* | 79 |
| 2022* | 97 |
| 2023* | 120 |
| 2024* | 147 |
| 2025* | 181 |

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- **Data Science has emerged as an important discipline**
- We need to tackle new challenges, and many more are yet to come
- We will be learn about modern day **pioneers** in our class moving forward
- However, most the basic principles relates to basics of Math, Statistics and Probability theories. **Tribute** to

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

**Zetabyte**

- 2010: 2
- 2011: 5
- 2012: 6.5
- 2013: 9
- 2014: 12.5
- 2015: 15.5
- 2016: 18
- 2017: 26
- 2018*: 33
- 2019*: 41
- 2020*: 64.2
- 2021*: 79
- 2022*: 97
- 2023*: 120
- 2024*: 147
- 2025*: 181
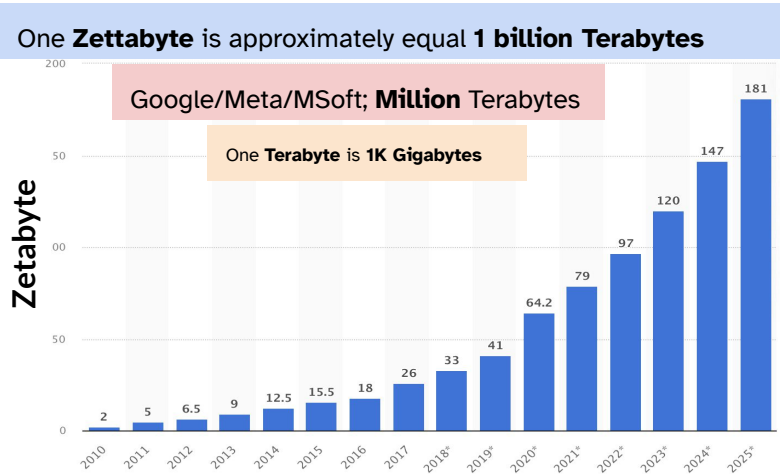
Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- Data Science has emerged as an important discipline
- **We need to tackle new challenges, and many more are yet to come**
- We will be learn about modern day **pioneers** in our class moving forward
- However, most the basic principles relates to basics of Math, Statistics and Probability theories. **Tribute** to
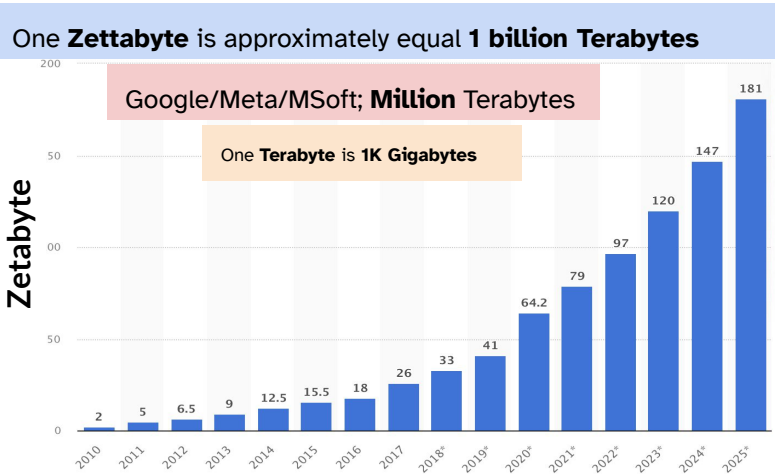
# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

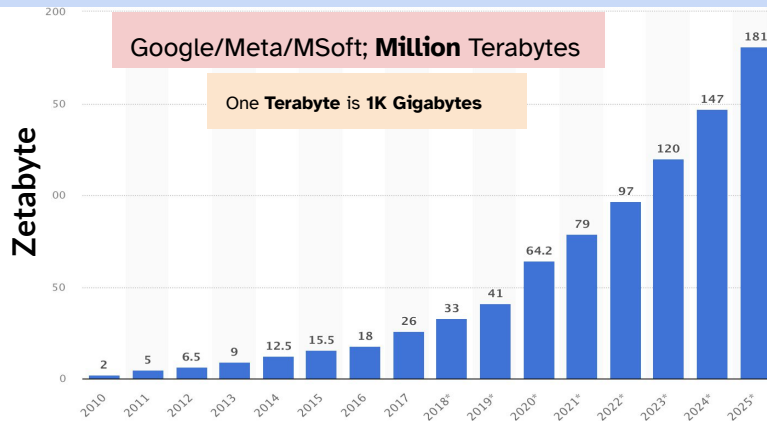One **Terabyte** is **1K Gigabytes**



Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025

- Data Science has emerged as an important discipline
- We need to tackle new challenges, and many more are yet to come
- **We will be learn about modern day pioneers in our class moving forward**
- However, most the basic principles relates to basics of Math, Statistics and Probability theories. **Tribute** to

# Data Growth and the Trend

One **Zettabyte** is approximately equal **1 billion Terabytes**

Google/Meta/MSoft; **Million** Terabytes

One **Terabyte** is **1K Gigabytes**

Zetabyte

| Year | Value |
|------|-------|
| 2010 | 2 |
| 2011 | 5 |
| 2012 | 6.5 |
| 2013 | 9 |
| 2014 | 12.5 |
| 2015 | 15.5 |
| 2016 | 18 |
| 2017 | 26 |
| 2018* | 33 |
| 2019* | 41 |
| 2020* | 64.2 |
| 2021* | 79 |
| 2022* | 97 |
| 2023* | 120 |
| 2024* | 147 |
| 2025* | 181 |

<u>Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025</u>

- **Data Science** has emerged as an important discipline
- We need to tackle new challenges, and many more are yet to come
- We will be learn about modern day **pioneers** in our class moving forward
- However, most the basic principles relates to basics of Math, Statistics and Probability theories. **Tribute** to
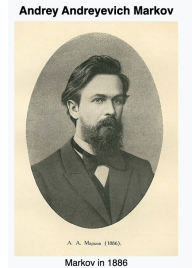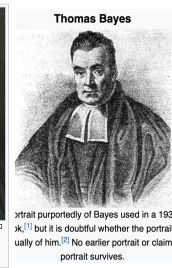


Bernoulli's *Ars Conjectandi* was the first work that dealt with probability theory as currently understood.

Carl Friedrich Gauss made major contributions to probabilistic methods leading to statistics;

Thomas Bayes
ortrait purportedly of Bayes used in a 1936 )k,[1] but it is doubtful whether the portrait is ually of him,[2] No earlier portrait or claimed portrait survives.

Gerolamo Cardano (16th century)

Andrey Andreyevich Markov
Markov in 1886

# QA