# **Topic 2**: Data/information modalities

© blog.tail.digital

# **August – December 2021**

*Juan Carlos Gómez*
PhD in Computer Science
jc.gomez@ugto.mx
http://jcgcarranza.wix.com/juancarlosgomez
Office 314

Campus Irapuato-Salamanca | División de Ingenierías

Structured

Images

Text

**Information Modality**

The **way**/form in which the information is **presented**.

Variable/Atribute

Structure

| Name | Sex | Age | Semester | Height |
|------|-----|-----|----------|--------|
| Juan | H | 19 | 3 | 1.70 |
| Ana | M | 20 | 5 | 1.65 |
| Maria | M | 18 | 1 | 1.56 |
| Pedro | H | 21 | 8 | 1.75 |

Observation/Instance

# **Structured Data**

The **structure** is separated from the content. That means, each observation is an **instance** of a previously defined structure

© J.C. Gomez

# **Structured data**

- Each **column** could represent a different type of **variable** (nominal, ordinal or interval)

- (Generally) **Easy** to enter, manipulate, store, read, query and process in a computer

# Structured data: processes

- **Statistical** description

  - Max $\quad \max\{x_i...n\}$

  - Min $\quad \min\{x_i...n\}$

  - Sum $\quad \sum_{i=1}^{n} x_i$

  - Mode $\quad \mathrm{mod}\{x_i...n\}$

  - Median $\quad median\{x_i...n\}$

  - Mean $\quad \frac{1}{n} * \sum_{i=1}^{n} x_i$

  - Variance $\quad \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{\mu})^2$

  - STD $\quad \sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i - \mu)^2}$

  - Unique $\quad unique\{x_i...n\}$

  $$P(G_i) = \left(\frac{cell\ count(G_i)}{n}\right)$$

  - Entropy $\quad \sum_{i=1}^{J} P(G_i)(-\ln(P(G_i)))$

  - Probability $\quad \sum_{i=1}^{J}(P(G_i))^2$

© keywordsuggest.org

© J.C. Gomez

# Structured data: processes

◎ **Plots**/charts (data visualization)

# **Structured data: task**

- ◉ **Variable associations**

Crime incidents in 2017 in Washington, D.C.



© bookdown.org

© J.C. Gomez

# Structured data: task

- Prediction

**Monthly crude oil prices (Jan 2016-Dec 2021)**
dollars per barrel



8

# Structured **data**: Issues

- If there are **many variables**, there will be difficulties to process them

- **Incomplete** data

- **Noise** (spurious data)

- **Outliers**

**Helen Palmer**

| From: | robert@rhxgroup.com.au |
|---|---|
| Sent: | Monday, 2 May 2011 4:35 PM |
| To: | helen@rhxgroup.com.au |
| Subject: | Business Documents |

Dear Helen

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Quisque pulvinar fermentum cursus. Sed quis magna sit amet nunc dignissim convallis. Duis mauris mi, ultricies eget varius ac, hendrerit ut nulla. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Aliquam sed lorem massa, non lobortis nulla. Curabitur orci dui, sollicitudin in tincidunt vitae, sollicitudin eget magna. Duis vitae tellus erat, sed tincidunt velit. Duis scelerisque tempus commodo. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Aliquam id lectus id justo vehicula rhoncus. Pellentesque sed interdum nunc. Ut fringilla semper semper. Nullam vel ante sit amet lectus varius tempus. Cras consequat tristique neque ac consectetur. Maecenas ante ipsum, bibendum ornare porttitor interdum, porttitor non ligula. Nam tincidunt, eros nec pretium euismod, diam augue tristique metus, sed accumsan ipsum arcu non sem.

Curabitur mollis est et nibh luctus tincidunt. Praesent ac turpis sed risus consectetur aliquam. Proin accumsan feugiat arcu, vel ornare elit dictum sed. Donec ipsum tortor, commodo a tempor nec, porttitor at nisl. Nam viverra euismod metus id rhoncus. Etiam iaculis aliquam volutpat. Nam eget sapien sit amet mi scelerisque egestas non eu elit. Suspendisse potenti. Maecenas sit amet ante risus. Maecenas commodo, leo et consectetur aliquam, velit mauris eleifend ante, a tristique justo massa eget justo. Curabitur a eros risus. Curabitur egestas diam velit, in facilisis quam. Nunc eleifend placerat suscipit. In blandit massa eu turpis eleifend ut laoreet purus sagittis. Quisque in turpis quis leo sodales aliquam in quis urna. Nulla facilisi. Donec eget urna in mauris ornare rutrum. In tellus risus, tempus eget feugiat a, hendrerit ac nibh.

I look forward to your response about the documents by the end of the week.

Kind regards
Robert

# Textual Data

A coherent set of written **signs** that **transmits** an informative message

# Textual data

⦿ **Raw** meaning: Set of symbols (letters/logograms).

⦿ **Semantic** meaning: Informative message content (what is the reader understanding from the text)

# 📌 Textual data

- Emails
- Social networks posts
- General documents
- Web pages
- Legal documents
- Medical reports
- Research papers
- Books
- …

**Helen Palmer**

| | |
|---|---|
| From: | robert@rhxgroup.com.au |
| Sent: | Monday, 2 May 2011 4:35 PM |
| To: | helen@rhxgroup.com.au |
| Subject: | Business Documents |

Dear Helen

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Quisque pulvinar fermentum cursus. Sed quis magna sit amet nunc dignissim convallis. Duis mauris mi, ultricies eget varius ac, hendrerit ut nulla. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Aliquam sed lorem massa, non lobortis nulla. Curabitur orci dui, sollicitudin in tincidunt vitae, sollicitudin eget magna. Duis vitae tellus erat, sed tincidunt velit. Duis scelerisque tempus commodo. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Aliquam id lectus id justo vehicula rhoncus. Pellentesque sed interdum nunc. Ut fringilla semper semper. Nullam vel ante sit amet lectus varius tempus. Cras consequat tristique neque ac consectetur. Maecenas ante ipsum, bibendum ornare porttitor interdum, porttitor non ligula. Nam tincidunt, eros nec pretium euismod, diam augue tristique metus, sed accumsan ipsum arcu non sem.

Curabitur mollis est et nibh luctus tincidunt. Praesent ac turpis sed risus consectetur aliquam. Proin accumsan feugiat arcu, vel ornare elit dictum sed. Donec ipsum tortor, commodo a tempor nec, porttitor at nisl. Nam viverra euismod metus id rhoncus. Etiam iaculis aliquam volutpat. Nam eget sapien sit amet mi scelerisque egestas non eu elit. Suspendisse potenti. Maecenas sit amet ante risus. Maecenas commodo, leo et consectetur aliquam, velit mauris eleifend ante, a tristique justo massa eget justo. Curabitur a eros risus. Curabitur egestas diam velit, in facilisis quam. Nunc eleifend placerat suscipit. In blandit massa eu turpis eleifend ut laoreet purus sagittis. Quisque in turpis quis leo sodales aliquam in quis urna. Nulla facilisi. Donec eget urna in mauris ornare rutrum. In tellus risus, tempus eget feugiat a, hendrerit ac nibh.

I look forward to your response about the documents by the end of the week.
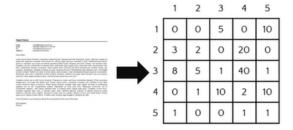
Kind regards
Robert

# Textual data: processes

### Word filtering



### Frequency analysis



### Transformation

# Textual data: tasks

© opentext.com

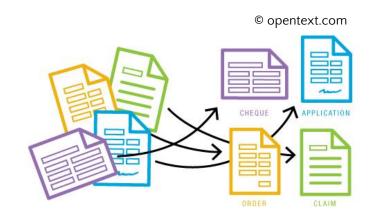**Classification**: Spam/phishing/legal email, news, advertises, web pages, social media posts, sentiment detection, etc.

'Infinity war', the last Avengers movie was great! → **Good**

'Infinity war' stars Robert Downey Jr as Iron Man. → **Neutral**

'Infinity war' was a horrible movie. → **Bad**

© J.C. Gomez

# 📌 Textual data: tasks

◉ **Analysis**: Plagiarism detection

© J.C. Gomez

# Textual **data**: issues

◉ **Large** number of words



5.3 million of articles in Wikipedia (English)

→

**3.2 billion words!**

© J.C. Gomez

# Textual **data**: issues

## ◉ Polysemy

All New 2018 **Jaguar** XF
The **jaguar** is an endangered species
Mac OS X 10.2 "**Jaguar**"

## ◉ Synonymy

The **big** dog has a **large** house
It makes me **sad** to see you
so **unhappy**

# Textual data: issues

- Language

© J.C. Gomez

# Different modalities, same observation

| Brand | Model | Year | HP | Cylinders |
|-------|-------|------|-----|-----------|
| BMW | 435i | 2018 | 248 | 4 |

Coupe and convertible variant of BMW's compact 3-series sedan, the 435i doesn't disappoint. It has a turbocharged engine of 248-hp 2.0-liter inline-four and can be equipped with either a six-speed manual or an eight-speed automatic and rear- or all-wheel drive.

**Each modality provides with specific (sometimes complementary) details about the observation**

© J.C. Gomez

# End topic 2

Next **topics**

- Analysis of structured data
- Analysis of text data
- Unsupervised learning

# Credits

Special thanks to all the people who made and released these awesome resources for free:
- ⦿Presentation template by SlidesCarnival