

Harvesting Digital Fields: Creating a Corpus of Web Scraped Data

Camila Lívio

Department of Research and Computational Data Management

University of Georgia

November 15, 2023

Table of Contents

- 1 Introduction
- 2 Objectives
- 3 Intensification
- 4 Two Studies
 - First Task
- 5 Future Research

Introduction

Three Wolf Moon T-Shirt

This shirt has changed my life! (...) with The Mountain Three Wolf Moon Short Sleeve Tee, the SHIRT now draws the eye. One young teen even shyly approached me, and instead of asking for a picture or an autograph, simply smiled conspiratorially and whispered, “Team Jacob, right? Me, too. He’s soooooo dreamy.” Yes he is, young lady. Yes. He. Is.

Tuscan Milk

Shipping was fine, and the product was not damaged in any way, but my husband and I (both of us have college degrees, mind you, his in Engineering) could not figure out how to assemble this. No instructions, no diagrams, not even a lousy cheap allen wrench. So basically we’re using it as a one gallon paper weight. I haven’t gotten any response from Tuscan. It earns two stars simply because it is heavy and does do a fair job of holding down the stack of newspapers awaiting recycling.

- 1 What kind of text are these?
- 2 What are the general characteristics of this text genre?

Vásquez, Camilla. “My life has changed forever!” Narrative identities in parodies of Amazon reviews.” *Narrative Inquiry* 27.2 (2017): 217-234.

Today's lecture:

- (1) Building a specialized corpus from web scraped data
- (2) Intensification
- (3) Some elements to consider when building a corpus

Intensifiers – A quick overview

- Speaker's tendency to innovate and engineer different ways to boost meaning
 - Scale up the quality or property of another word
 - Go beyond semantic gradability as they increase the expressive strength and subjectivity in discourse
 - **'Really cool', 'Absolutely amazing', 'A proper disaster'**
- What are other intensifiers do you hear these days (in any language you speak/work on)?

Semantic gradability & Subjectivity

‘Absolutely amazing movie’ vs. ‘an amazing movie’

Their illocutionary force varies: statement and an expressive act

Semantic gradability & Subjectivity

cont.

Intensifiers have also been studied in light of social membership

Martos (2010): *super* has ascended, socially and stylistically, to be perceived as a middle and upper-class female usage in Madrid (2010: 202)

Bradac et al. (1995): ‘female register’ or ‘register of low power’. The use of intensifiers, hedges, and tag questions as linguistic forms that are commonly associated with women

Inventory

- ❶ Adverbs (*muito/muy, bem/bien*, ‘very’)
- ❷ Multi-word expressions (*um monte de/montón de*, ‘a lot of’)
- ❸ Affixes (*Pt. belíssimo/trilegal*, *Sp. riquísimo, superguay, superinteresante*)
- ❹ Absolute adjectives (*Perfeito/perfecto* ‘perfect’)
- ❺ Pronouns (*Trae unos Raleighs, ¡córrele!*)¹
- ❻ Certain verbal constructions (*¡Vaya coche tan moderno!*, ‘What a modern car!’)²
- ❼ Que + adjective (*Que lindo!*, ‘So + Adjectives’)
- ❽ Comparative phrases (*Sp. Más feliz que una lombriz*, *Pt. Mais quebrado que arroz de terceira*, ‘Happier than/Cooler than’, ‘As cool as...’)

“intensity is signaled by a large and miscellaneous class of devices, ranging from the most peripheral of prosodic variations to the most central categories of the grammar” (Labov 1985a: 48).

Two Studies

(1) Intensifiers in the RCor = Reviews Corpus

(2) Social Evaluation of Intensifiers

Why Intensifiers in Online Reviews?

- Intensifiers are very common in evaluative language
- Language that expresses an attitude towards a person, a situation or any entity, being both subjective and situated in a societal value-system (Hunston 2010)
- Evaluative language = subjective
- The scalar nature of evaluative language
- The connection between evaluative language and its scalar nature offers a strong motivation for considering customers' reviews as a relevant linguistic source for the study of intensifiers
- Online reviews as a text genre

Building the RCor

- R Core Team (2020)
- Web scraping
- Specialized corpus (Brezina and McEnery 2020: 17-19)
- Brazil and Mexico are the two countries in Latin America with the highest number of visits to Amazon.com
- RCor is balanced across product types
- Filter reviews written in the target language
- Organization of the data in columns and rows
- Data quality check
- Half a million words



Run, or he's going to tell us about
again!

R

What does the RCor look like?

Figure: RCor – MS

review_id	review_text	review_star	review_date	review_format
1	Perfecto	5.0 de 5 estrellas	Revisado en México el 3 de enero de 2017	Color: Azul/azulante del producto: empaquetado está
2	Excelente	5.0 de 5 estrellas	Revisado en México el 26 de abril de 2019	Cámara verificada
3	¡Justo increíble el color! Y super bien el ajuste	5.0 de 5 estrellas	Revisado en México el 18 de noviembre de 2020	Color: Teal Gold-V07781Talla: MCompra verif
4	Muy Rápido	5.0 de 5 estrellas	Revisado en México el 28 de diciembre de 2020	N/A
5	¡La mejor "cámara" que he comprado! ¡La amor	5.0 de 5 estrellas	Revisado en México el 16 de enero de 2021	Color: Negro (dark black)Cámara verificada
6	Sopero. De animales a Dioses. Una breva histo...	4.0 de 5 estrellas	Revisado en México el 16 de diciembre de 2021	Cámara verificada
7	¡Buena cámara!	5.0 de 5 estrellas	Revisado en México el 21 de diciembre de 2020	Cámara verificada
8	¡Uno de lo mejor!	5.0 de 5 estrellas	Revisado en el Reino Unido el 7 de diciembre	Cámara verificada
9	Más de lo que esperaba	5.0 de 5 estrellas	Revisado en México el 10 de agosto de 2018	N/A
10	¡Mi favorito	5.0 de 5 estrellas	Revisado en México el 1 de agosto de 2021	Cámara verificada
11	Producto bonito en la foto de cámara	5.0 de 5 estrellas	Revisado en México el 19 de diciembre de 2020	Cámara verificada
12	Muy útil	5.0 de 5 estrellas	Revisado en México el 12 de julio de 2020	Denominación: 1036Nombre del disco: Serris
13	Me encanta	5.0 de 5 estrellas	Revisado en México el 3 de agosto de 2021	Cámara verificada
14	¡AMAZON! ¡Excelente!	5.0 de 5 estrellas	Revisado en México el 22 de febrero de 2021	Cámara verificada
15	Muy fácil de manejar, instalar y usar	5.0 de 5 estrellas	Revisado en México el 7 de febrero de 2021	Cámara verificada
16	Muy preciso y calidad	5.0 de 5 estrellas	Revisado en México el 11 de diciembre de 2019	Color: Silver-BlackCámara verificada
17	Buen servicio	5.0 de 5 estrellas	Revisado en México el 2 de marzo de 2016	Cámara verificada
18	Buena cámara	5.0 de 5 estrellas	Revisado en México el 23 de septiembre de 2...	Cámara verificada
19	Buena producto	5.0 de 5 estrellas	Revisado en México el 24 de junio de 2020	Cámara verificada
20	Cámara	5.0 de 5 estrellas	Revisado en México el 16 de marzo de 2021	Cámara verificada

- Review title
- Text of the review
- Star rating
- Review date
- Review format

Scraping the Data

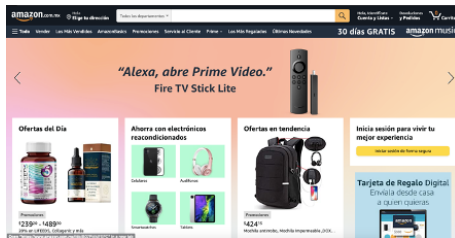
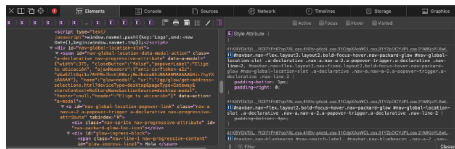


Figure: Screenshot of an Amazon page



Building a corpus via web scraping

- Data scraping: a technique in which a computer program extracts data from human-readable output coming from another program.
 - The scraper was programmed to collect random reviews from 50 pages for each product
 - Top-seller of each category (*Mais vendidos*, ‘best sellers’)
 - **Benefits**: amount of data, and the ability to gather and organize the data set in a way that facilitates looking for patterns through the use of corpus linguistics methods.
 - RCor is balanced across product categories
- What would be some alternative ways to balance this type of corpus?

Ethical Considerations

- **Drawbacks:** Data set is *not* publicly available so that the research is compliant with ethical considerations, such as personal data and privacy, and copyright and intellectual property
- Do not collect any personally identifiable information
- Use the data only for the purpose for which it was collected
- Use aggregated data rather than individual reviews (helps with maintaining anonymity)
- Check the website's '**robots.txt**'
- Look for information about writing a Data Management Plan (DMP)
- GitHub

Research Questions – Task 1: Intensifiers in Online Reviews

- ① How do intensifiers vary across product categories?
- ② Are there any associations that can be made between the use of intensifiers in positive and negative Amazon reviews?
- ③ What can the study of intensifiers in this environment teach us about computer-mediated communication?

The analysis is divided in four main parts

- (1) Prefixation, Suffixation and Lexical Intensifiers
- (2) Intensifiers by Product Categories
- (3) Intensifiers by Product Ratings
- (4) Intensifiers and the construction of a reliable online persona

Intensifiers by Product Category

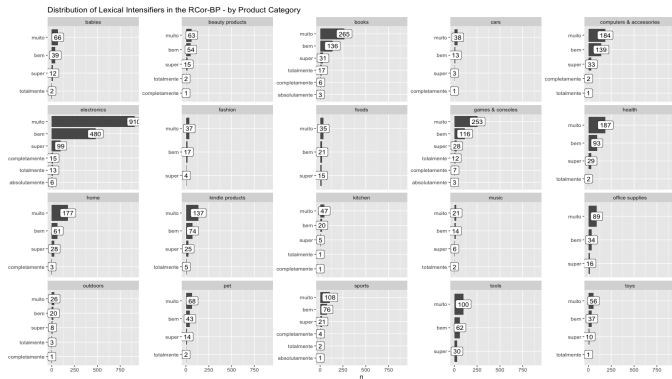
Question

How is the use of intensifiers different, if at all, across different product categories?

- Some product categories show a higher frequency of intensifiers, as well as a more varied inventory
- Felbermayr and Nanopoulos (2016): where there is language related to play and fun, the emotional spectrum tends to be larger (i.e., a wider set of words that describe emotions)
- BP: books, beauty, computers and electronics, games, home, kindle, kitchen, and music == categories with over 3 different intensifying forms
- MS: amazon devices, baby, books, beauty, home, musical instruments, toys == categories with over 3 different intensifying forms

Intensifiers by Product Category

Figure: Intensifiers by Product Category



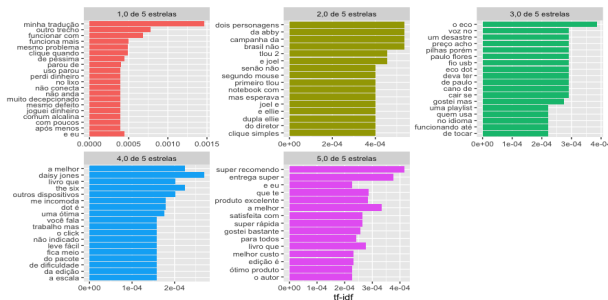
Intensifiers by Product Rating

Question

Does rating matter?

- Plotted the correlation between bigrams extracted from positive (4-5 stars) and negative (1-2 stars)
- Created two tf-idfs
- Combinations that appear more than five times
- Intensifiers are preferred when reviews are positive
- When reviews are negative, users seem to focus on the description of the product (and less on their experience)
- Alternative way to visualize this connection: distribution of bigrams by product rating

Figure: Bigrams by Rating in the RCor-BP



Conclusions of Task 1

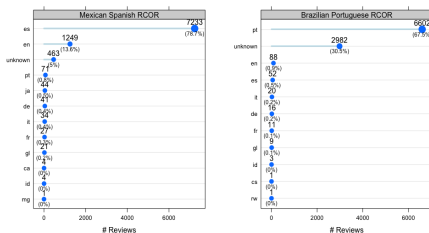
- Online reviews can be considered a text genre
- Intensifiers are a crucial feature of such genre
- I connect the compulsory use of intensifiers in online reviews to the user's need to be relevant, and show credibility (i.e., expressing their commitment)
- Canonical intensifiers are most frequent, a fact that I attribute to the user's construction of a reliable persona, trustworthiness, and expertise (Mackiewicz 2008)
- Digital self-expression to state 'entrepreneurial agency' > Weninger and Li (2002) mention the existence of a 'microcelebrity persona' > somebody who is willing to bear the cost of broadcasting to achieve some kind of recognition
- Hu et al. (2009): the overconfident user is the one who is more likely to bear the costs of broadcasting > explain the overwhelming number of positive reviews

- Patterns of spelling variation and paralinguistic resources: emojis and caps lock (McCulloch 2019, 2020)
- Quantitative patterns of stacking in the Romance languages (Scheffler et al. 2023)
- Can tense influence the distribution of intensifiers?
- The relevance of intensifiers in the timely subject of fake news (Grieve and Woodfield 2023)

Super obrigada!
camila.emidio25@uga.edu

Languages in the RCor prior to Filtering

List of languages found in the RCor corpora



Star Rating Distribution - Note the j-shaped distribution (Hu et al. 2009)

