# Text Analytics & Business Application

Text Analytics in E-commerce

Qinglai He

Department of Operations and Information Management

Wisconsin School of Business

WISCONSIN SCHOOL OF BUSINESS

# IC5 Discussions

- Why using pre-trained word embedding as text representation could not give us better classification results than BoW?

  - The effectiveness of word embeddings and BoW depends on various factors, such as the data and complexity of the task. So, it is always beneficial to experiment with different approaches for a particular task.

  - The performance of pre-trained word embeddings also depends on other factors, such as overlap of words between data the pre-trained model, corpus used for training, training architectural, and its dimensionality.

  - Pre-trained model aims to capture more general meaning. In contrast, BoW is generated based on data in a specific domain, therefore it is more flexible for different tasks.

# IC5 Discussions

- What is the pros and cons of using word embeddings?
  - Pros:
    - It captures more semantic relationships between words. Therefore, it is good for tasks where the meaning of words is crucial.
    - It can capture more contextual information.
    - It provides dense and low-dimensional representations.
  - Cons:
    - Fixed vocabulary: The vocabulary of the pre-trained model are from training corpus. It may not cover all words in a specific domain or dataset.
    - May not capture the domain-specific nuances, limiting the effectiveness for certain tasks.
    - Fine-tuning the pre-trained model is usually required for a specific domain.
    - We are unable to interpret the embeddings.
    - Still couldn't capture the word order in your given sentence.

# IC5 Discussions

- Available pre-trained model offered by Gensim
  - Read more here https://radimrehurek.com/gensim/models/word2vec.html

## Pretrained models

Gensim comes with several already pre-trained models, in the Gensim-data repository:

```
>>> import gensim.downloader
>>> # Show all available models in gensim-data
>>> print(list(gensim.downloader.info()['models'].keys()))
['fasttext-wiki-news-subwords-300',
 'conceptnet-numberbatch-17-06-300',
 'word2vec-ruscorpora-300',
 'word2vec-google-news-300',
 'glove-wiki-gigaword-50',
 'glove-wiki-gigaword-100',
 'glove-wiki-gigaword-200',
 'glove-wiki-gigaword-300',
 'glove-twitter-25',
 'glove-twitter-50',
 'glove-twitter-100',
 'glove-twitter-200',
 '__testing_word2vec-matrix-synopsis']
```

# Outline of Today's Class

- E-Commerce applications

- E-Commerce catalog

- Sentiment analysis

  - Review analysis

  - Aspect-Level sentiment analysis

- Case study

# Text Analytics Applications in Ecommerce

# Recommendation System

# Delighting Customers with Personalized Discounts

# Fighting Cyberfraud



**PROMOTION ABUSE**
Individual unduly benefits from a coupon several times.

**AFFILIATE FRAUD**
Faudster generates false activity to receive commissions from an affiliate program.

**FRIENDLY FRAUD**
Merchant receives a chargeback because the cardholder denied receiving the order, yet the goods were actually received.

**RESHIPPING FRAUD**
Fraudster recruits a person to package and re-ship merchandise purchased with stolen credit cards.

**IDENTITY THEFT**
Fraudster buys or steals sensitive personal information to buy goods.

**ACCOUNT TAKEOVER**
Fraudster is using another person's account information to obtain products and services.

# E-commerce Catalog

# E-commerce Catalog

- A product catalog is a database of the products the enterprise deals with, containing product descriptions and images.

- Better product descriptions with relevant information help customers choose the right product through the catalog.

- Product information can also help with search and recommendations.

# Building an E-Commerce Catalog

- Building an informative catalog is one of the primary problems in e-commerce. It can be split into several subproblems:

  - Attribute extraction

  - Product categorization and taxonomy creation

  - Product enrichment

  - Product deduplication and matching

# 1. Attribute Extraction

- Attributes define product properties in e-commerce and provide a complete overview of the product for the customer.

- A rich set of attributes improves clicks and click-through rates, which influences product sales.

- Obtaining attributes directly from sellers for all products is difficult, and the quality of attributes must be **consistent** for the customer to have the correct and relevant information.



REFINE BY

CLOTHING
OUTERWEAR (4)
POLO SHIRTS (4)
SHORTS (0)

COLOUR
☐ BLACK (2)
☐ BLUE (3)
☐ GREEN (2)
☐ GREY (2)
LIGHT (0)
☐ PURPLE (1)
☐ WHITE (5)

SIZE
☐ S (4)
☐ M (4)
☐ L (4)
☐ XL (4)
☐ XXL (1)

CLEAR ALL

**ADIDAS**
FP ENGINEERED STRIPE POLO
$69.00

# Obtaining Product Attributes

- Traditionally, e-commerce websites employed manual labeling or crowdsourcing techniques to obtain product attributes, which is expensive and not scalable with an increase in the volume of products.

- Machine learning techniques are used to obtain product attributes.

- There are two types of attribute extraction algorithms:
  - Direct (will talk about information extraction in Week 9)
  - Derived (Indirect, e.g., classification task)

# 2. Product Categorization and Taxonomy

- Product categorization is a process of dividing products into groups.

- These groups can be defined based on similarity
  - e.g., products of the same brand or products of the same type can be grouped together. Generally, e-commerce has pre-defined broad categories of products, such as electronics, personal care products, and foods.

# Product Categorization and Taxonomy

- A good taxonomy and properly linked products can be critical for business operations.

- This categorization process is typically manual to start at small scale.

- At scale, this categorization is typically posed as a classification task where the algorithm takes information from a variety of sources and applies the classification technique to solve it.

# 3. Product Enrichment

- To improve search and recommendations in e-commerce platforms, it's important to gather richer product information.

- Short and long titles, product images, and product descriptions are potential sources of this information.

- However, this information is often <span style="color:red">incorrect</span> or <span style="color:red">incomplete</span>, which can hamper the faceted search in e-commerce platforms.

# Product Enrichment: Misleading Titles

- Misleading titles can hamper faceted search in e-commerce platforms and negatively affect click-through and conversion rates.



This text is too complicated even for a human to parse and make sense of, let alone a machine.
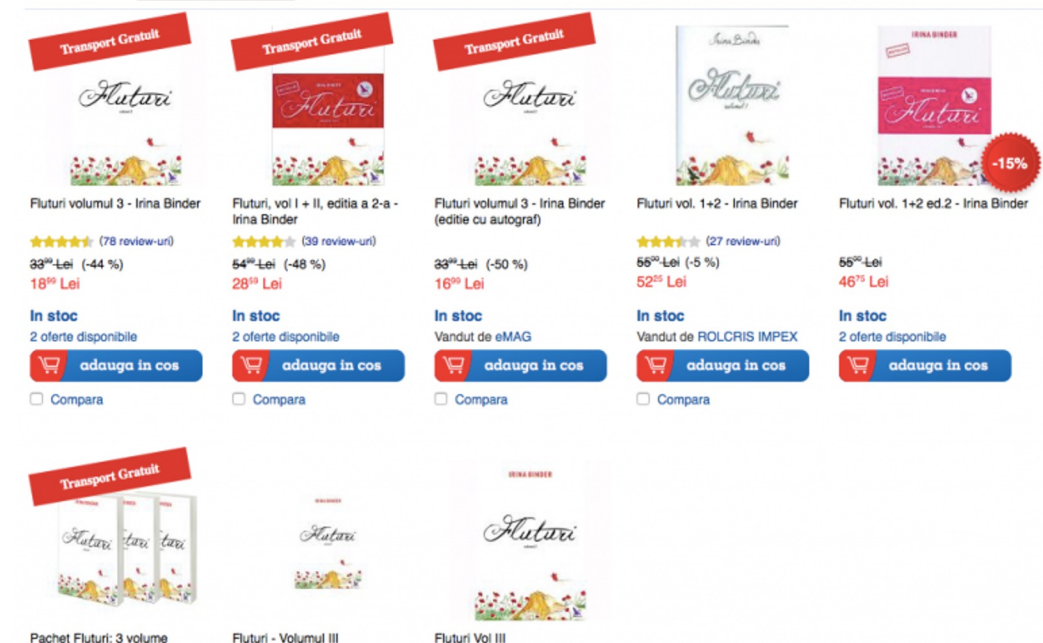
# Improving Product Titles

- Filtering out misleading tokens

  - <span style="color:red">Direct string matching</span> can be used to improve product titles.

  - Tokens that are not part of the product's attribute values should be filtered out to avoid misleading search results.

- Using pre-defined templates

  - <span style="color:red">Pre-defined templates</span> can help maintain consistency across product titles.

  - A template can be built using attributes from the taxonomy tree.

# 4. Product Deduplication and Matching

- Product duplication in e-commerce
  - Third-party sellers often add products to e-commerce platforms with different names and terminology.
  - This can result in the same product being listed with multiple titles and product images, leading to confusion for customers.

- Ways to handle product duplication
  - Attribute match
  - Title match
  - Image match

# Attribute Match

- Matching attribute values

  - Maximum overlap of attributes indicates strong product matching.

- Using string matching

  - String matching can be used to match attribute values.

  - Two strings can be matched using exact character match or string similarity metrics that take into account slight spelling mistakes, abbreviations, etc.

# Title Match

- One product can often have multiple title variants.

- A simple method could be to compare **bigrams** and **trigrams** among these.

- It's also possible to generate title-level features (such as counts of common bigrams and trigrams) and then calculate the Euclidean distance between them

# Image Match

- Different techniques like pixel-to-pixel match, feature map matching, and advanced image-matching techniques can be used for image matching.



**Dress**

**Description**
Walker
Neckline Round neckline
Pattern Striped
config_sku M9121C3CM-Q11
Fit Normal
Length Knew length
Sleeve length Short sleeves
M9121C3CM-Q11
-25%

**Dress 'Walker-H'**

**Description**
Crew neck
Quilted hem/ edge
Straight hem
straight cut
Chest pocket
striped
All-over pattern
Article no: 430237391999
Material

Elastic/ stretch
Viertelarm
Short/ Mini
regular fit

**WALKER - Everyday dress**

**Description**
Combined design

# Sentiment Analysis

# Why Sentiment Analysis?

- **Movie**: is this review positive or negative?

- **Products**: what do people think about the new iPhone?

- **Public** sentiment: how is consumer confidence? Is despair increasing?

- **Politics**: what do people think about this candidate or issue?

- **Prediction**: predict election outcomes or market trends from sentiment

# 1. Review Analysis



**Day Pack**
Top Topics

| | reviews | stars | rating |
|---|---|---|---|
| attractive | | | |
| durable | | | |
| comfortable | | | |
| roomy | | | |
| features | | | |
| lightweight | | | |
| computer | | | |
| travel | | | |

- They capture direct feedback from customers about products, and they can directly affect the sales of the products.

- **Brand monitoring**: reveal key insights about how your brand, product, or company is viewed by your customers and stakeholders

- **Product improvement**: companies can use them to further improve the customer experience.

# 2. Aspect-Level Sentiment Analysis

- An aspect is a collection of words that indicates certain properties or characteristics of the product.

- Aspects can also include anything related to the supply, presentation, delivery, return, quality, etc.

- Supervised and unsupervised techniques:
  - Supervised techniques involve using labeled data and predefined features.
  - Unsupervised techniques involve clustering and topic modeling.

# Challenges of Sentiment Analyses

- Tone

- Polarity

- Sarcasm

- …

# Case Study & Tools

# Existing Tools for Sentiment Analyses

- TextBlob -  an easy way to calculate sentiment polarity. TextBlob is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.

- Vader Sentiment Analysis

- Entity Sentiment Analysis

# Examples - TextBlob

```
text_1 = "The material feels amazing and it seems durable"

TextBlob(text_1).sentiment.polarity
```

=> 0.6

```
text_2 = "I have purchased two of these dresses in the past and was
 excited to order one from Amazon. Unfortunately quality is very lo
w. Arrived with stitching already unraveling."

TextBlob(text_2).sentiment.polarity
```

=> -0.09375

# Examples with Google Colab

**Take 10 minutes break…**