# How we use NLP at Barion

Budapest ML 2022

**Zoltán Balogh, PhD**
Senior Data Scientist
zoltan.balogh@barion.com

# CONTENTS

- About **Barion** & Data Monetization
- Problem statement
- Challanges of training set preparation
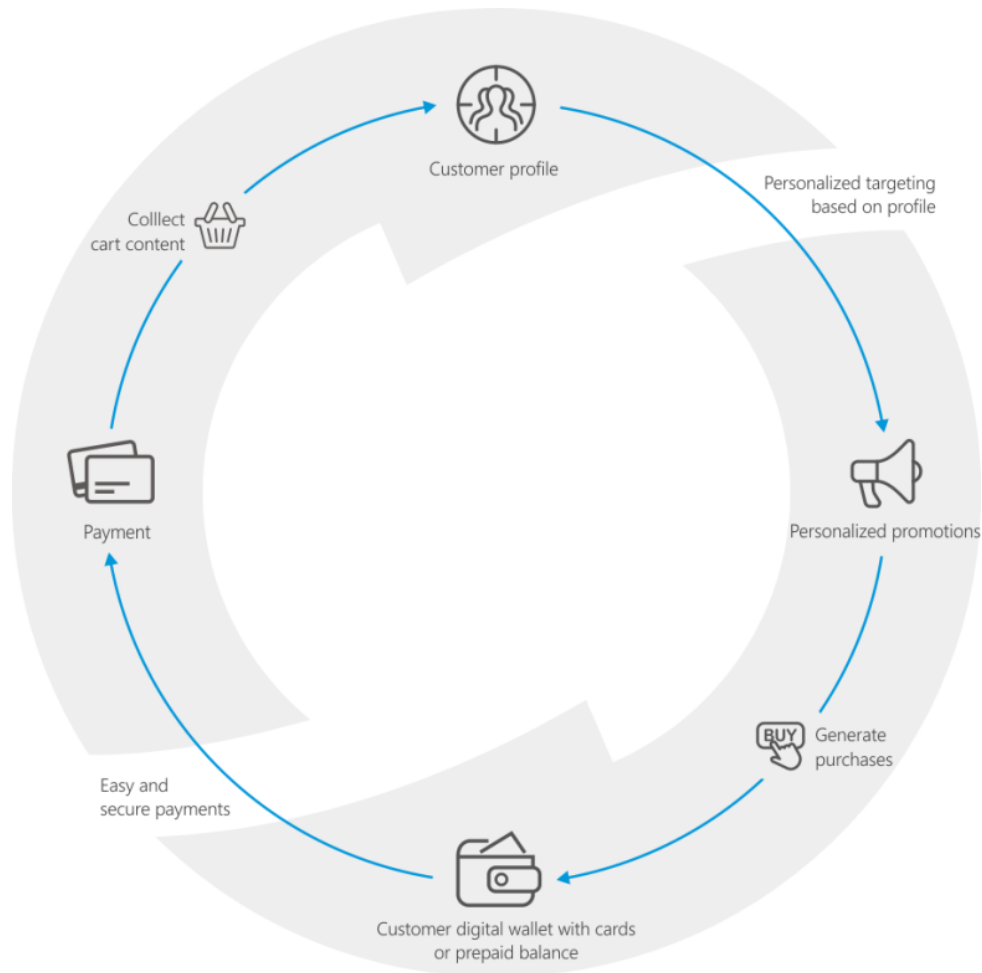- Model building & training
- Used softwares

# About Barion

- Licensed as an electronic money issuer
- With **Barion Smart Gateway** customers can easily and safely pay in more than 13 000 online shops mainly in CEE
- In exchange for **lower payment fees**, merchants can choose to **share consented data** of their **customers**, which is collected and stored in our data lake

# Data Monetization

o **Barion Pixel** is a Javascript snippet built into the **merchant's website**

o With the customers' prior consent, the **details** of their **shopping behaviour** is sent to Barion

o The collected **data** is then **transformed** into profiles

o The created profiles are utilized in different advertising campaigns to enhance their targeting options

o When the customers visit webpages, **personalized promotions** are **shown** to them



Colllect cart content

Customer profile

Personalized targeting based on profile

Personalized promotions

Generate purchases

Customer digital wallet with cards or prepaid balance

Easy and secure payments

Payment

# Properties of Collected Data

- BarionPixel collects different events to capture the customers' shopping behaviour from more than 2000 merchants
- In this presentation we will be focusing on the contents of their shopping basket
- The items of the customers' basket are individually sent to and evaluated by the neural network
- The raw input of the neural network is the **category** and the **name** of the product or service provided by the shop
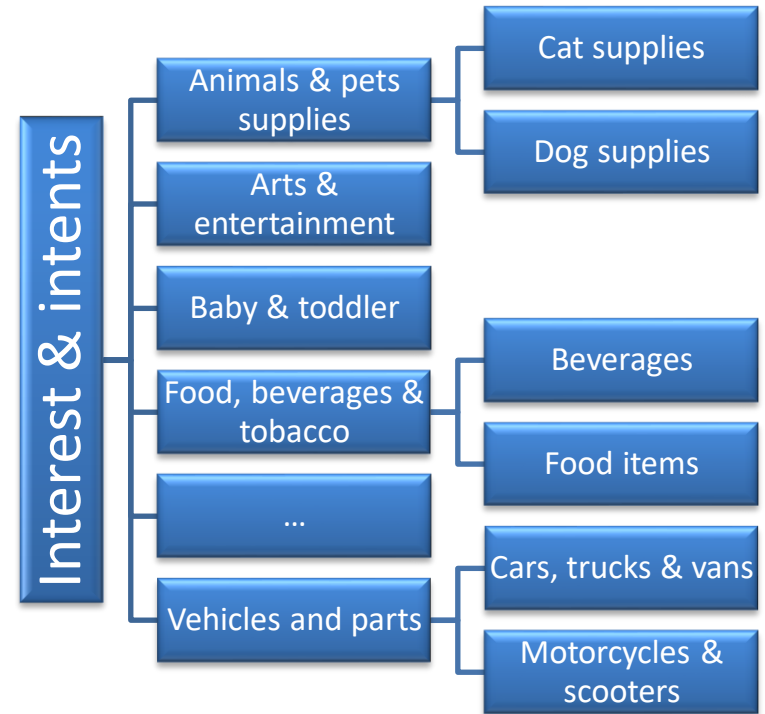
# Requirements of the Output

o Simple keyword search is not accurate and resource efficient
o A quick, robust and resource efficient model is needed that achieves high accuracy

o The output marketing segment should be produced from the model output with adjustable granularity
- o Ranging
  - o From **getting acquainted** with the **product** (visited the product only once)
  - o To the customers who **bought** the **product** or **have a long history of interest** for it
- o Accuracy/Confidence property

# Segmentation Bases

o Marketing segmentation can be based on
  o Geolocation
  o Technographic
  o Demographic
  o **Interest and Intent**

o Models
  o **Neural network**
  o Regex

# Challanges of NLP training set preparation

# Peek into the raw input data

o **Audi A6 autó izzó** — Disproportionate keywords within sample
o Újszülött pelenkák - Pop-in - Gyártók - Zöld Úton - Mosható
o Férfi MTB kerékpár - MTB kerékpárok - Kerékpárok - 5 - Bicaj
o AlphaOne DZ Inteligentní hodinky, bílé
o **Ledvance Planon Plus 30W 2700-6500K 595x595mm felületre szerelhető LED panel távirányítóval** — Noisy inputs
o **Női** felsők, **klassz** pólók hölgyeknek **kedvező áron**
o Line Sugar Effect Gel Silcare - Nail4U
o **Canon PG545XL** — Keyword not present
o **Asus P8H61-M LX2 1155 alaplap+ CPU hűtő** — Keywords might belong to a different category
o Street Surfing Ripper Roller - Bloody Gold
o **;Névre szóló baba- Angyalka (Új) Karácsonyi limitált kiadás** — Category not clear
o **R15 Gyűrű**
o **CN-HG 95 lánc** — Ambiguous

# Process of the preparation of the input

**Training set creation**
- Fetch stratified samples of products from relevant shops
- Keep only the keywords of the segmentations

**Conversion**
- Keeping the top most important keywords
- Vectorizing (100D)
- Vector size limit
- Embedding weights

**Encoding/splitting**
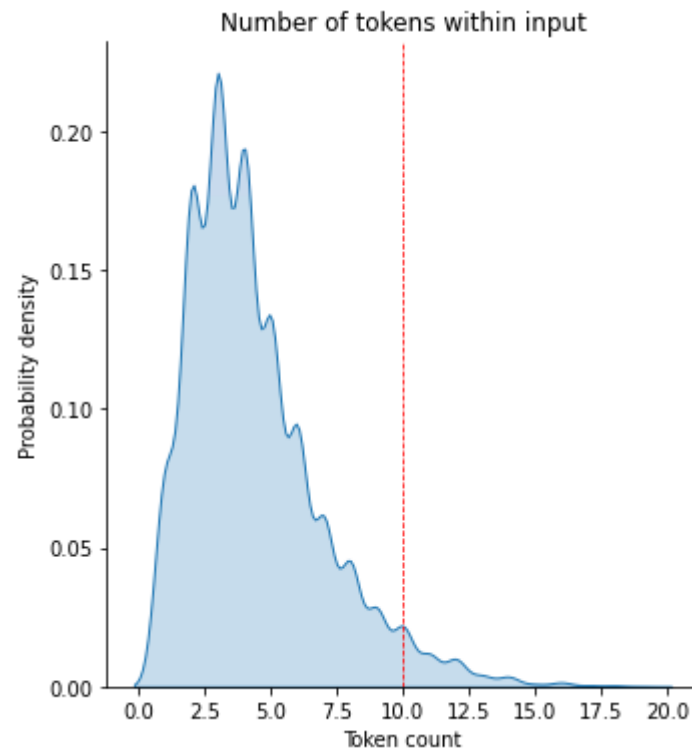- One-hot encode the output
- Train/test split

# Sampling

o Not all the categories have the same amount of quality entries
    o Stratified by the category of product/service
o The categories with excessive number of items need to be undersampled
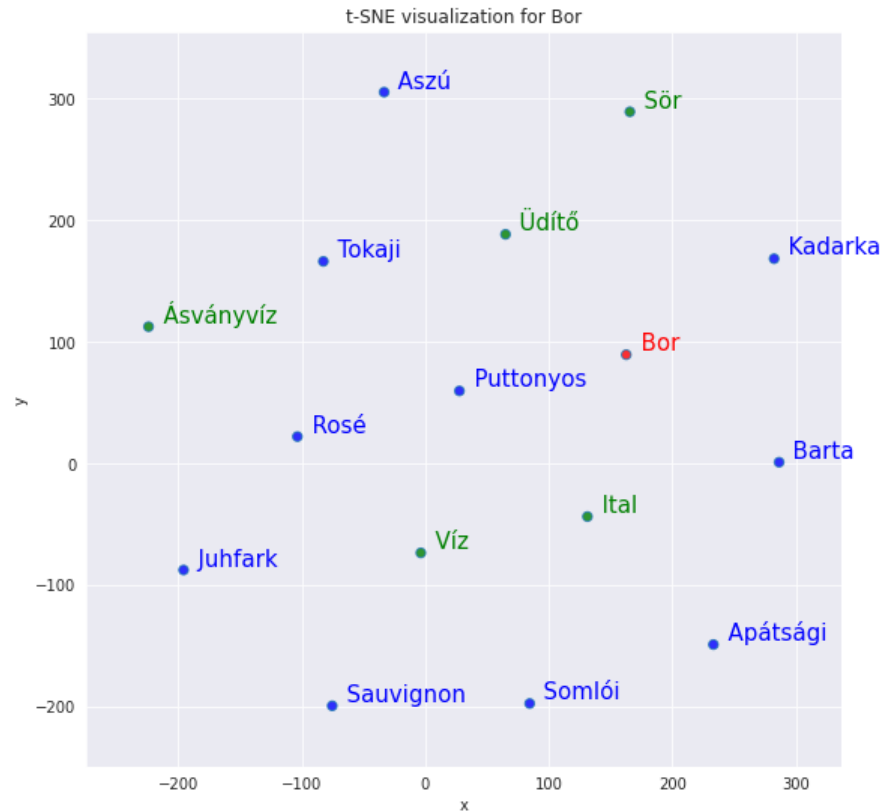    o Ensuring the include all the unique keywords relevant to the category

| Manually select shops for each category | Stratified fetch from database | Select all entries with unique keywords | Manually check the records in sample |
|---|---|---|---|

# Token Cleansing

o Only the first 10 tokens within the input text
  is kept (as only 3.02% input entries are longer
  that 10 tokens)
o Stopword removal (tailored to the Hungarian
  market)
  o Colors (grey, yellow, green etc)
  o Measurement units (pcs, liter, xl, month)
  o Conjunctions (and, or)
  o Others (new, compatible, super, import,
    export, action, premium etc.)



Number of tokens within input

# Vectorization

o Vectorizing using gensim.Word2Vec
o Converts all tokens to a 100 dimensional vector space
o Ensures that words with similar meaning are closer to each other in space



t-SNE visualization for Bor

# Model architecture

o Requirements
- o The model should output only 1 category (no overlaps)
    - o Binary crossentropy loss function
- o Needs to capture the word order in phrases effectively

o Tried different models of 2/3 layers
- o LSTM with heavy regularization and dropout
- o CNN
- o LSTM bidirectional

**Lábápoló  csiszoló  gyűrű**

Jewellery: apparel and accessories

Tools: hardware

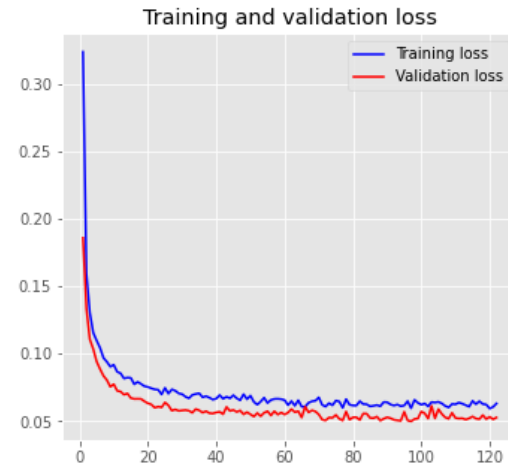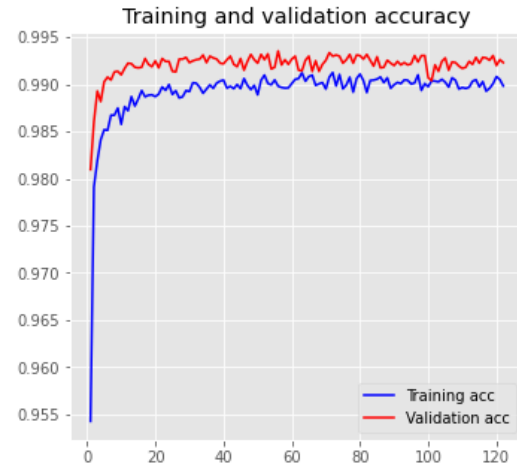Toe polishing tool: Health and beauty

**Baba szappan** vs **szappan baba**

Soap: Health and beauty

Toy: Toys and games

# Training & Output

o The output of the model and other features are saved into database linked to the user
o The output is not necessarily sales-ready
  o The database needs to be queried for a specific segment
  o The output can be customized further on different business demands



Training and validation accuracy



Training and validation loss

# Different Architectures

o Sequential model
o First embedding layer with the weights of the pre-trained 100 dimensional embedding matrix

o Dense layer with sigmoid activation
o Adam optimizer, binary crossentropy loss function, optimize for accuracy

|  | LSTM | CNN |
|---|---|---|
| Number of layers | 2-3 | |
| Regularization | L1, dropout | |
| Filters | 256, 128, (64/32) | |
| Extra | | Global max pooling |

The 2-layer LSTM performed the best

# Result

- In recent a campaign for a **well-known multinational electronics company** promoting their child-care product-line, we provided our **„Baby and toddler"** category
- The targeted audiences in the given campaign were provided by **Google** and **Barion**
- The Barion segments outputted by the neural network gained **~25%** better performance, than the audiences provided by Google

# Softwares, we use

**barion**

Thank you for your attention!

**Zoltán Balogh, PhD**

Senior Data Scientist

zoltan.balogh@barion.com

www.barion.com