

# An application of Naive Bayes on Amazon's products

Arthur Bonetti

*Master in Finance - Financial Entrepreneurship and Data Science*

*University of Lausanne*

Lausanne, Switzerland

arthur.bonetti@unil.ch

**Abstract**—Amazon is a vast marketplace where it can be difficult to find the perfect product that meets our need. This project proposes the implementation of a Naive Bayes classification algorithm on a small subset of Amazon's product to advise better products to an user. The process is optimized by looping on small sections of the database and performing pre-processing algorithm.

**Index Terms**—Naive Bayes, Amazon, database, scraper, optimization, classification

## I. INTRODUCTION

One of the yearly struggle that everyone has already experienced once is the trouble of finding a good gift for a relative or friend. With a lot of marketplaces available such as Amazon, Aliexpress, Zalando or Digitech and million of products to choose from, one can easily spend hours searching and lose himself in the vast expanse of products available.

Algorithms used by these giants re-sellers are already at the cutting edge of the technology and companies like Amazon spend 42 million a year on technology and content expenses [1]. However, the drawback for the user, with these algorithms, is the constant monitoring of the activity on the website. Every second an user spend searching for the perfect gift is tracked, from how long he looked at a product to his search activities history. The aim of the company is clear, it is to make the user spend more time and money on their website by advising complimentary products among others. The kind act of looking for a gift for a relatives can quickly become a psychological battle between algorithms and self control.

Furthermore, the shopping activity does not end once the user left the company's website. The activity is tracked and is used to advertise products on other websites. Someone who just wanted to buy baby socks for his friend's newborn can easily be advised baby products for a few weeks.

## II. RESEARCH QUESTION

### A. Problem

A current problem for some people when searching for a product online is over information. According to a paper published in 2017 [2], with the internet providing massive amount of heterogeneous information, people may perceive this medium as challenging. Information overload can be a

consequence of the difficulty to select and evaluate relevant information. The paper was conducted in the context of online news exposure. In a safe way, this can be extended to the selection of a gift while online shopping.

### B. Objective

The objective of this paper is to create a program for people looking into buying a gift in an easy way, without overwhelming them with information. The aim is to propose products based on a like or dislike test and to recommend new gifts based on a Naive Bayes classification algorithm. Naive Bayes is a classification method that is often used for text classification, but variants for arbitrary data exist. Naive Bayes relies on the Bayes theorem and the assumption that the data is independent. The end result of our program is to redirect the user directly on the page of the product he has chosen so that he can buy the gift he has selected without being tracked by the company's website.

### C. Scope

Considering the size of the product's database and the difference of information available among online shopping companies, this project's scope lies on a self made database from a subset of products from Amazon. The database contains 950 products from eight different categories. The scope can easily be extended to larger database and different online shopping website.

## III. METHODOLOGY

### A. Database

The first step into the project is to create an adequate database to serve as an example and a basis for testing the program capabilities. The initial idea was to use the Application User Interface (API) from Aliexpress to retrieve data from products. However, this process required to own a seller's account and was not adapted for the scope of this project. Therefore, we decided to use a tool available on amzscout [3]. The tool amzscout Product Database allow to fetch product from different categories and obtain useful information on the product features. The product database is limited to a few searches on the free account, but this can easily be overridden by the use of temporary emails and the creation of a few free accounts. We then collected the different fetched data and combined them into a single .csv file.

### B. Scraping

The second step is to scrape images for each product from Amazon's website. Indeed, having a picture of the product is an essential element to help the user decide on the likeness of a gift. Pictures usually are larger in term of memory size than common data, therefore they are not stored in the database. Because we need to show the product's picture to the user, we decided to scrap the main picture from the product's URL. This is done using the `imgscrape` library from Pypi.

### C. Visual Interface

The project needs a way to show the information to the user. To this end we use a convenient library named Tkinter that allows us to create windows and add pictures, label and buttons with Python. The window is updated every time the user interacts with the interface. To keep a trace of the user's likeness of a product, we update a subset of our database every time the user interacts with the program.

### D. Naive Bayes Classification

To advise better products to the user, we use a Naive Bayes classification algorithm. At first, we don't have information about the user's preference. Therefore, we select ten products at random from the categories that are selected. Then, the user decides for each product whether he likes it or not. If the user likes the product, it is classified as 1 and if he dislikes it, it is classified as 0. After ten iterations, we compute a probability score for each product on the whole data set using the classification algorithm and fetch the ten products obtaining the score closest or equal to one. The next step is to show the new products to the user and continue the process. Every ten iterations, we compute again the classification probability with the new observations and thus provide a better classification every time.

## IV. IMPLEMENTATION

The program for this project was programmed using Python 3.8.3 and requires `chromedriver` and the following packages to run:

- Pandas - version 1.2.3
- Numpy - version 1.19.2
- Tkinter - version 8.6
- Selenium - version 4.0
- Pillow - version 8.2.0
- Requests - version 2.25.1
- Mkl-random - version 1.1.1
- Sklearn - version 0.0
- Matplotlib - version 3.3.4
- Imgscrape - version 0.1.1

### A. Pre-processing algorithm

Pre-processing the data is a necessary step for the correct working of the program. We implement the cleaning, scraping and saving of the data in the `scrape_n_save.py` file. First, we import the raw data from the `data_prog.csv` file to a pandas dataframe. We then drop the column that are not useful for

our program to reduce the amount of data. We also drop the columns containing empty values that might causes the classification algorithm to complain later on. After that we create two functions. The first function `scrape_img(lst_var)` calls the `scrape` function from the `imgscrape` package, gives it a list of URLs and return the result. The second function `iter_image(df, df3)` iterates over the database, fetch the URL for each product from the database by bulk of ten, save them in a smaller dataframe and call the `scrape_img(lst_var)` function. When this function is called the computer will use `chromedriver` to open a chrome tab, go to each product's URL and save the product's main image URL in a variable. When the program has fetched the pictures' URL from the ten products, it adds them to the dataframe under a new column named `Image` and save the dataframe as a csv.

At first, we wanted to scrape the image of the product during the execution of the program so that we were not stocking the pictures' URL for each product. However, it appeared that scraping the pictures was a time consuming process as the computer has to open a chrome tab for every request. We considered the use of parallelism for this process, but it seems that it is not quite a good idea to open ten google chrome tabs at the same time and the computer won't allow it anyway.

We fetch the picture by bulk of tens to offer a guarantee in case the scraper encounters an error. Indeed, a product can be deleted from Amazon or the computer can go to sleep during the process. Saving the results every ten iteration gives the possibility to take back the process for where it crashed when an unexpected error occurred instead of having to run it from the beginning again. For information, during the first run of this file and with raw data, the program encountered two URLs leading to missing product, went to sleep twice and was once stopped by an accidental touch of the keyboard by the programmer.

Once the pictures' URL are collected into the respective .csv files, we import the csv to dataframes and concatenate them together to have the final database that will be used for the program. The database is saved in the `datav_prog_final.csv` file.

### B. Keeping data of interest

For this project, we put some effort to keep the size of the data used during the program to the minimum. We have eight product's categories in our database and we know that some categories will interest some user and not others. Therefore, at the beginning of the main program, we ask the user for the categories of interest. Because it is time consuming and not efficient to write the name of the categories, we simply ask for numbers. The user can input the categories he wants to keep. Each number is attributed to a category. For example one particular user could write "12378". We then use a loop and a `str()` function to see if the input contains numbers.

The numbers are then converted to create a new dataframe that keeps only the categories of interest that will be used for the program. The new dataframe created is named df2. The algorithm used to analyse the input will keep only the numbers. Therefore, if an user gives as an input "aiàèç2kl45", the algorithm will detect 2, 4 and 5 and keeps the affiliated category. If an user does not give any input, we consider that he wants to keep all categories.

After that step, we select ten rows at random from df2 and save them in a new dataframe named df3. This new dataset will serve as a basis to display the information in the Tkinter window and compute the classification score.

	Name	Price	Min Price	Net	FBA Fees	LQS	Category	Sellers			
	827 AstroAI Air Comj	25.49	25.49	16.15	9.34	78	Tools & Home Improvement	1			
	842 Filtrrete 10x20x1	80.06	57.25	59.08	20.98	80	Tools & Home Improvement	4			
	817 Rust-Oleum 249	4.79	4.79	-0.85	5.64	84	Tools & Home Improvement	10			
	570 PowerA Play & C	14.88	11.99	9.48	5.4	87	Video Games	4			
	556 Nintendo Blue/ f	66.99	66.99	57.42	9.57	84	Video Games	23			
	392 Avildove Lingerie	15.99	15.99	9.06	6.93	90	Clothing, Shoes & Jewelry	3			
	144 Insect Lore Butte	36.99	36.99	26.34	10.65	84	Toys & Games	2			
	659 WORX WG896 1	93.53	86.05	64.12	29.41	74	Patio, Lawn & Garden	12			
	32 Zmodo SD-H108	30.6	30.6	21.08	9.52	68	Electronics   Computers   Acces	2			
	851 Filtrrete UR10, fil	101.94	101.94	76.12	25.82	86	Tools & Home Improvement	3			
Rank	Est. Sales	Est. Revenue	Reviews	Col Rating	Weight	Colors	Sizes	URL	Image	Score	ML score
20	94560	25.49	4396	4.4	2.2	2	1	https://www.https://imaj	1	0	0
82231	73880	5412640	30	4.9	0.9	1	1	https://www.https://imaj	0	0	0
51	34228	441061.1	34214	4.7	1	127	1	https://www.https://imaj	0	0	0
496	2243	33866.99	5912	4.7	0.48	6	1	https://www.https://imaj	0	0	0
45	9751	653219.49	63026	4.8	0.37	6	1	https://www.https://imaj	0	0	0
434	8501	135930.99	24268	4.5	0.13	26	7	https://www.https://imaj	0	0	0
10	26939	996473.61	13214	4.7	1.3	1	1	https://www.https://imaj	0	0	0
405	8332	779291.96	5518	4.4	16	1	1	https://www.https://imaj	0	0	0
9	44292	1355335.2	1859	4.3	0.8	1	1	https://www.https://imaj	0	0	0
33	32914	3355253.2	17384	4.8	4.34	1	26	https://www.https://imaj	0	0	0

Fig. 1. Dataframe df3 example

### C. Tkinter interface

We use Tkinter to show the information about a product to the user. Our dataframe contains eighteen columns and most of the information is not necessary for the user to decide of his likeness to a product. The information that we consider essential for the user is :

- The picture of the product
- The description of the product
- The price of the product

We consider these three items as an object, therefore we create a class named Model(object) that keeps the above mentioned information in a single object. In the same class, we also initialize a counter that allows us to keep track of the number of products that has been shown to the user. We get the values from the first row of the dataframe df3 and create a Model(object) named count\_object. During the whole process of the program, we will use this count\_object and update it with the values of the current product.

When the program is launched, we initialize the Tkinter window for the first time. We get the picture's URL from our object and use the request function to save the picture in a variable. After that, the picture is resized to a 400x400 format and is added to the window. We also get the description and the price of the product and create two labels under the picture.

We then create four buttons.

- Button "Yes" which allows the user to say that he likes a product and save the value 1 in the dataframe
- Button "No" which allows the user to say that he does not like a product
- Button "Stop" which allows the user to stop the program
- Button "Buy" which redirects the user to the page of the product on the merchant website

Once these elements are placed in the Tkinter window, we show them to the user.

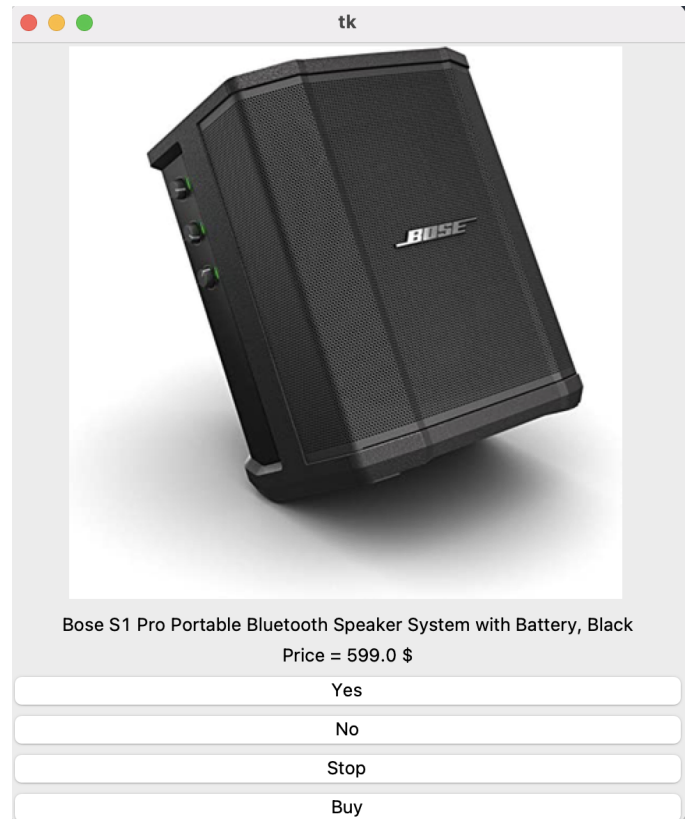


Fig. 2. Tkinter window example

Now, we create a function update\_tk() that updates the window to show new information every time the user click on Yes or No. The function updates the counter of the current object and updates the values of the object to be the values of the next row in the dataframe df3. Once this is done, we update the picture and the labels accordingly and we keep the current layout. This process is repeated until the user reaches the last (tenth) product of the dataframe df3. As a reminder, every time that the user clicks on "Yes", the value 1 is saved in a column named "Score". The default value of the "Score" column is 0 so that when the user clicks on "No", the program does not have to update the dataframe.

### D. Naive Bayes Classification

Once the user has gone over the ten first product, we predict the probabilities of a product from the dataframe

df2 to be classified as 1. To this end, we create a function named `ml_score(df3, df2)`. In this function, we first define which columns we wish to consider for the X values and we collect these values from the dataframe `df3` and save them to a variable X. We then define our regressor y which is the values of the column "Score". The score's value is 1 if the user liked the product and 0 otherwise. Therefore, we are trying to predict which other product could have a score of 1 based on the X observations. The values we are trying to predict are in the dataframe `df2`, which as a reminder is the dataframe that contains the information on all products from categories that the user decided to keep. From this dataframe, we select the same columns as in X and create our variable `X_test`. We then use a `preprocessing.MinMaxScaler()` to normalize the data and fit our data in a Naive Bayes Classification algorithm from the `sklearn` package. The last step is to predict the probabilities of the product to belong to the class 1. The values are stored in the dataframe to the "MI\_score" column.

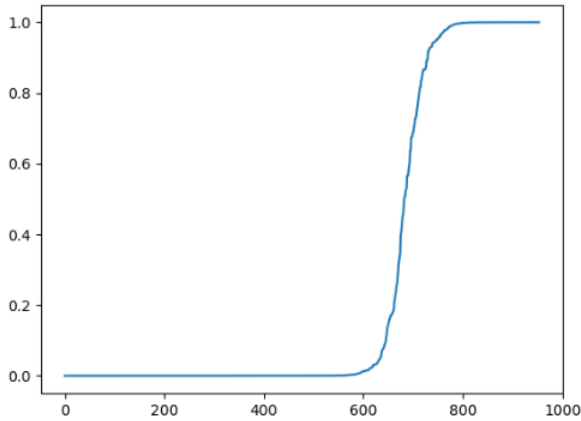


Fig. 3. Naive Bayes Classification example with 3 Yes out of 10

#### E. Updating the dataframe df3

Now that we know the predicted score for our whole dataframe, we select the ten products that got the biggest predicted score with the Naive Bayes Classification algorithm and add them to the small dataframe `df3`. Because we use this dataframe a lot throughout the program, we decided that it was best to keep track of it by creating a class named `Model_df(object)`. We store the dataframe as a `Model_object` and update it every time we get the ten new best results with the classification algorithm.

Because we do not want to fetch products that we already showed to our user, we drop the product's row from our dataframe `df2` that we add to the dataframe `df3`. Another advantage of that process is that we get rid of data duplicates and therefore save time on computation.

#### F. Finding the dream gift

For every iteration, we collect more user's preference and can compute the Naive Bayes Classification algorithm with more observations. The result of this is a more precise classification every time the user's gives his opinion on a product by clicking on "Yes". The program shows products that are more likely to interest the user and he will, more often than not, find the perfect gift that he came looking for. Once he has found that gift, the user simply has to click on the buy button to be redirected to the web page of the product. This is simply done by using `chromedriver` and giving the product's URL as an argument.

#### G. Performance

The program for this project does not require a lot of processing power and does not need a lot of time to process the computation. Indeed, the program is written in a way that each loop iterates for a maximum of ten times. The computation in the loop are quite basics and does not require a lot of resources from the computer. The overall improvement in performance from the first version of the project was done by pre-processing the data and scrape the product's picture URL in advance so that the user has a lag free experience. The project can easily be replicated on a bigger dataset without causing trouble. However, we have to keep in mind that the Naive Bayes Classification algorithm predicts values for the dataframe `df2`. It can cause some slow-down in the execution of the program if the size of this dataframe is big enough. In that case, we would advise to perform the classification algorithm on a smaller dataframe. Because we only need to fetch ten new products every time, an enormous dataframe is not needed.

### V. MAINTENANCE AND UPDATE

At the moment, the program requires the manual input of the category's names and the manual selection of the column's names of the dataset for the Naive Bayes Classification algorithm. Therefore, an interested third party should update these variables to meet his needs before launching the program.

The codebase for this program is located on my personal git repository: `LaPetiteBiche/prog_final_project`. Anyone interested can create a pull request to update the code. The last git pull is a working version of the program with a lot of comments added to help any third parties understand the code. The README is also updated and gives the necessary information for every person interested in this program.

### VI. RESULTS

The result of this project is a working program that helps the user find the perfect gift. To illustrate the result of the project, we will follow the steps of an hypothetical user. The user launches the program and has first to decide of the categories he wants to keep. Our user wants to find a product for his nephew who likes video games and sport, therefore he selects the categories 1. Toys and Games, 3. Sport, Outdoor and

Fitness, 4. Video Games, 8. Electronics. He puts as an input for the program "1348", then he accidentally pushes "jbkn" before submitting. The program recognizes that the user wants the categories 1, 3, 4 and 8 and creates a dataframe df2 containing only the product from these categories.

```
Welcome to the gift selector !
-----
We have 8 gift categories, please select the categories that you are interested in.
-----
1. Toys and Games, 2. Clothing, Shoes and Jewelry, 3. Sports, outdoor and Fitness, 4. Video Games
, 5. Patio, Lawn and Garden, 6. Musical Instrument, 7. Tools and Home, 8. Electronics
Please type the numbers of the categories (eg:12378) 1348
```

Fig. 4. Welcome Message and Input example

The resulting df2 is a 18 columns and 463 rows dataframe. At that stage, the user is shown a product and has to decide if he likes it or not. The first product he gets is a T-shirt. He dislikes it, clicks on "No" and continues the process. On the 10 first products, he only likes one, which is a Minecraft Lego set.



Fig. 5. Tkinter window, Minecraft Lego set

After 10 products, the program computes the probability score for the 453 remaining products to have a score of 1 with the Naive Bayes Classification algorithm. Because the user only liked one product so far, the algorithm classifies most of the product to belong to the group 0.

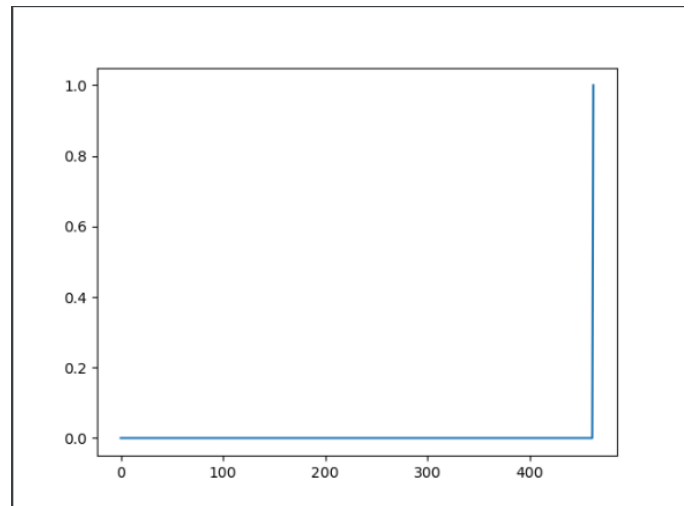


Fig. 6. Predicted Classification Score after 10 products

The program fetches the higher score and adds them to the products dataframe df3 to show them to the user. The user continues to likes and dislikes products and this time found 5 products that he liked. The program executes again the classification algorithm and chooses the 10 next best score to add to the df3.

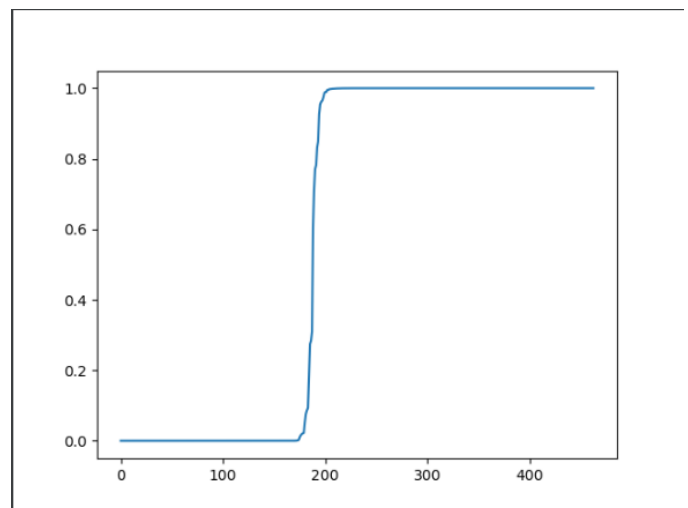


Fig. 7. Predicted Classification Score after 20 products

This time, with more observations, the program can create a new estimation for the predicted scores. Because the user liked more products, the program can with more confidence classifies more products to belong to the category with a score of 1.

For the 21th product, the program seems to have understood that the user likes products similar to the Minecraft Lego dataset from the beginning and proposes another Minecraft Lego set.





Fig. 8. Tkinter Window, another Minecraft Lego set

The user continues to use the program and after a few products, a new classification score is computed. The classification score is computed in the background and it is impossible for the user to notice when this occurs. For him, the experience is lag free.

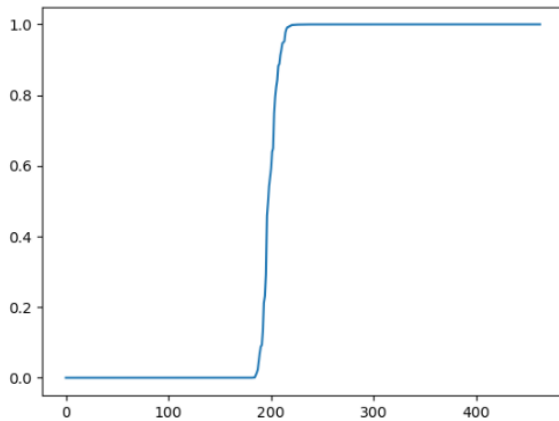


Fig. 9. Predicted Classification Score after 30 products

The Classification algorithm takes into account the new observations and computes new predicted values. Again, it fetches the predicted values with the highest score and proposes new products to the user. This time, the user is shown another Lego set with dinosaur Fossils. He fell in love with the product and decides to offer it as a gift to his nephew.

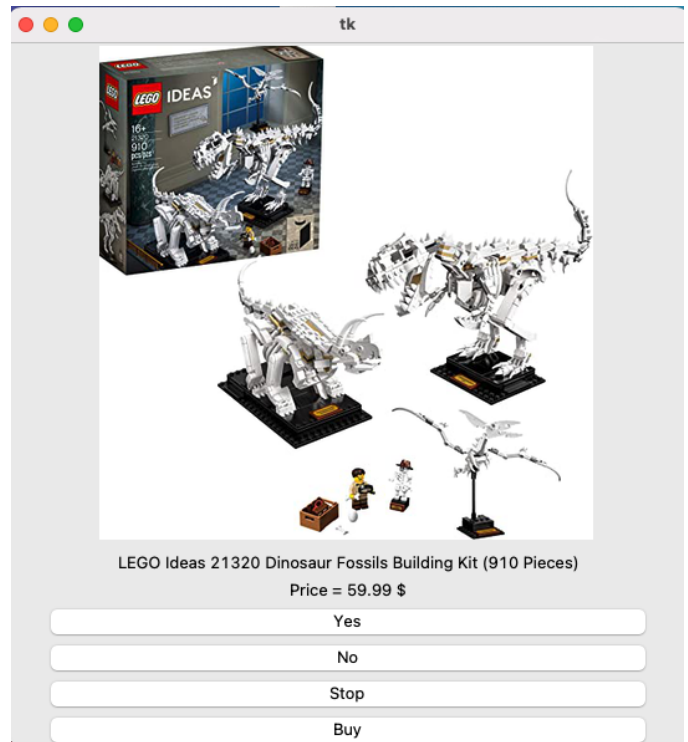


Fig. 10. Tkinter Window, Dinosaur Fossils Lego set

He clicks on "Buy" and he is directly redirected to the product's page on Amazon. He thanks the creator of the program for this experience and go on with his life. The full process only took 2 minutes, far less than what he would have taken using directly the website and a lot easier too. The user is impressed by the lag free experience and the accuracy of the products proposed.

## VII. CONCLUSION

This project resulted in a working version of the wanted result. The process is easier than using a regular medium to look for a gift. Therefore, we can affirm that this project has accomplished the goal of overcoming information overload. Furthermore, the program allows the user to bypass the websites and the tracking resulting from their use. On top of that, the Naive Bayes Classification algorithm seems to be greatly adapted to this project and offers meaningful results.

The program can certainly be improved. One way to improve the program can be to implement more general functions so that it can work with multiple databases without having to input manually the column names for the classification algorithm or the categories' names. It can also be interesting to implement other classification algorithms such as Gaussian Process classification and compare the results between the two algorithms.

The improvement in performance was done using the brain rather than parallel implementation. Indeed, the project was thought in a way that every loop is optimized so that only

relevant information is taken into account. The loops only execute ten times at a time with easy computation. The biggest optimization on performance was done by pre-processing the data and scraping the pictures' URL in advance so that this process had not to be done by the final user.

Overall, this project has attained his objectives and can still be improved visually and be generalize to more databases. At this stage, this is more a prototype than an end product that could be commercialized to the public. Nonetheless, it showed that it was possible with a database and a Naive Bayes Classification algorithm to show meaningful products based on an user's preferences.

## VIII. APPENDIX: CODE EXAMPLES AND ALGORITHMS

### A. Welcome\_message.py

```
1 # Welcome message function that return the input
2 def welcome():
3     print("Welcome to the gift selector !")
4     print("-----")
5     print("We have 8 gift categories, please select
6     the categories that you are interested in.")
7     print("-----")
8     print("1. Toys and Games, 2. Clothing, Shoes and
9     Jewelry, 3. Sports, outdoor and Fitness, 4.
10    Video Games")
11    print(", 5.Patio, Lawn and Garden, 6. Musical
12    Instrument, 7.Tools and Home, 8. Electronics")
13    category = input("Please type the numbers of the
14    categories (eg:12378)")
15    return(category)
```

### B. Scrape\_n\_save.py

```
1 import pandas as pd
2 import numpy as np
3 from imgscrape import imgscrape
4
5 # Create empty dataframe and give path to
6 # chromedriver
7 df3 = pd.DataFrame()
8 any_var = "/Users/arthur/Desktop/chromedriver"
9
10 # Import raw data
11 df = pd.read_csv('data/data_prog.csv')
12
13 # Cleaning - Deleting info we don't need
14 df = df.drop("Brand",axis=1)
15 df = df.drop("Available From",axis=1)
16 df = df.drop("Seller",axis=1)
17 df = df.drop("Shipping Weight",axis=1)
18 df = df.drop("Size",axis=1)
19 df = df.drop("Oversize",axis=1)
20 df = df.drop("ASIN",axis=1)
21 df = df.dropna()
22
23 # Scrap function
24 def scrape_img(lst_var):
25     var = imgscrape(*lst_var, path=any_var)
26     return(var)
27
28 # Iterating over dataframe 10 rows at a time and
29 # save image into df
30 def iter_image(df, df3) :
31     x = 10
32     while x < df.shape[0] :
33         for i in range(x-10,x):
```

```
32         df3 = df3.append((df.iloc[[i]]))
33         lst_var_df = df3['URL']
34         lst_var = lst_var_df.values.tolist()
35         y = scrape_img(lst_var)
36         m = np.asarray(y)
37         df3["Image"] = m
38         df3.to_csv(f'data/data_prog_{x}.csv')
39         x+=10
40         df3 = pd.DataFrame()
41
42 # Call function
43 iter_image(df, df3)
44
45 # Create variables
46 x=10
47 liste = []
48
49 # Create a list with all the csv
50 while x < 890 :
51     df = pd.read_csv(f'data/data_prog_{x}.csv')
52     liste.append(df)
53     x+=10
54
55 # Concat the csv and do some additional cleaning ->
56 # save the end result
57 result = pd.concat(liste)
58 result.to_csv(f'data/data_prog_final.csv')
59 df = pd.read_csv(f'data/data_prog_final.csv')
60 df = df.drop(df.columns[0],axis=1)
61 df = df.drop(df.columns[0],axis=1)
62 df.to_csv(f'data/data_prog_final.csv')
```

### C. plot.py

```
1 from matplotlib import pyplot as plt
2
3 # plot function
4 def plot_mlscore(ml_score) :
5     ml_score.sort()
6     plt.plot(ml_score)
7     plt.show()
```

### D. main.py

```
1 import pandas as pd
2 import numpy as np
3 import welcome_message
4 import random
5 import tkinter
6 from selenium import webdriver
7 from PIL import Image
8 import PIL.ImageTk
9 import requests
10 from sklearn import preprocessing
11 from sklearn.naive_bayes import *
12 from plot import *
13
14 # Create variables and create categories with
15 # categories name from dataset
16 categories = [
17     "Toys & Games", "Clothing, Shoes & Jewelry", "
18     Sports | Outdoors & Fitness",
19     "Video Games", "Patio, Lawn & Garden", "Musical
20     Instruments",
21     "Tools & Home Improvement", "Electronics |
22     Computers | Accessories"]
23
24 categories_final = []
25 df2 = pd.DataFrame()
26 df3 = pd.DataFrame()
27
28 # Path to chromedriver
29 any_var = "/Users/arthurbonetti/Desktop/chromedriver
30 "
```

```

26 # Import pre processed data from scv
27 df = pd.read_csv('data/data_prog_final.csv')
28
29 # Welcome message + input categories
30 categories_wanted = str(welcome_message.welcome())
31
32 # Save wanted categories in a list
33 for i in range(0, 9):
34     if str(i) in categories_wanted:
35         categories_final.append(categories[i - 1])
36     else:
37         continue
38
39 if categories_final == [] :
40     for i in range(0, 9):
41         categories_final.append(categories[i-1])
42
43 # Create new dataframe df2 to keep only categories
44 # of interest
45 for i in categories_final:
46     temp_df = df[df['Category'] == i]
47     df2 = df2.append(temp_df)
48
49 # Reset index
50 df2 = df2.reset_index(drop=True)
51
52 # Select 10 random elements from df2 and save into
53 # df3
54 for i in range(10):
55     x = random.randint(0, len(df2)-1)
56     df3 = df3.append((df2.iloc[[x]]))
57
58 # Drop element in df2 so that we don't choose them
59 # again
60 for i in df3[[]]:
61     df2.drop(i)
62
63 # New column df3 for ML
64 df3['Score'] = 0
65 df3['Ml_score'] = 0
66
67 # Keep info we have to show
68 var = df3['Image'].values.tolist()
69 var2 = df3['Name'].values.tolist()
70 var3 = df3['Price'].values.tolist()
71
72 # Create Model with count, image, price, description
73 class Model(object):
74     def __init__(self, img, descr, price):
75         self.currentObject = 0
76         self.img = img
77         self.descr = descr
78         self.price = price
79
80 # Initialise object with first element
81 count_object = Model(var[0], var2[0], var3[0])
82
83 # Object to stock dataframe
84 class Model_df(object):
85     def __init__(self, data_f, data_f2):
86         self.data_f = data_f
87         self.data_f2 = data_f2
88
89 # Create object with current df3
90 my_df = Model_df(df3, df2)
91
92 # Function to compute Naive Bayes Classification
93 # algorithm
94 def ml_score(df3, df2) :
95     # Columns name we use as X variables
96     dfml = df3[['Price', 'Net', 'FBA Fees', 'LQS', 'Sellers', 'Rank', 'Est. Sales', 'Est. Revenue', 'Reviews Count', 'Rating', 'Weight', 'Colors', 'Sizes']]
97
98     X = dfml.values
99     # y is the score value
100     y = df3['Score'].values
101     # Same variable as before for the test dataset
102     X_test = df2[['Price', 'Net', 'FBA Fees', 'LQS', 'Sellers', 'Rank', 'Est. Sales', 'Est. Revenue', 'Reviews Count', 'Rating', 'Weight', 'Colors', 'Sizes']]
103     # Normalize the data
104     min_max_scaler = preprocessing.MinMaxScaler()
105     min_max_scaler.fit(X_test)
106     X = min_max_scaler.transform(X)
107     X_test = min_max_scaler.transform(X_test)
108     # initialize the Naive Bayes and fit
109     clf = GaussianNB()
110     model = clf.fit(X=X, y=y)
111     # Compute the probabilities for score
112     y_proba = model.predict_proba(X_test)[: ,1]
113     print(y_proba)
114     # Save probabilities in new column and plot
115     m = np.asarray(y_proba)
116     df2["Ml_score"] = m
117     plot_mlscore(m)
118
119 # Function to update df3
120 def update_df3(df3, df2, var, var2, var3) :
121     print(df2)
122     # Add the 10 elements with biggest predicted
123     # values from df2 and remove from df2
124     for i in range(10) :
125         x = df2['Ml_score'].idxmax()
126         df3 = df3.append(df2.loc[x])
127         df2 = df2.drop(x)
128
129     # Add the new elements to the list used to show
130     # the product
131     temp = df3['Image'].tail(10).values.tolist()
132     temp2 = df3['Name'].tail(10).values.tolist()
133     temp3 = df3['Price'].tail(10).values.tolist()
134
135     for i in range(10) :
136         var.append(temp[i])
137         var2.append(temp2[i])
138         var3.append(temp3[i])
139
140     # Reset index and fill empty score values with 0
141     df3 = df3.reset_index(drop=True)
142     df3['Score'] = df3['Score'].fillna(0)
143     # Update object with new df3
144     my_df.data_f = df3
145     my_df.data_f2 = df2
146
147 # Tkinter
148 # Update window with new infos
149 def update_tk(var_img, var2, var3, count):
150     # Update the product count
151     count.currentObject += 1
152     # Update infos
153     count.img = var_img[count.currentObject]
154     count.descr = var2[count.currentObject]
155     count.price = var3[count.currentObject]
156     # Update window
157     im = Image.open(requests.get(count.img, stream=True).raw)
158     im = im.resize((400, 400), Image.ANTIALIAS)
159     photo = PIL.ImageTk.PhotoImage(image=im)
160     canvas.create_image(0, 0, image=photo, anchor=tkinter.NW)
161     name.config(text=count.descr)
162     price.config(text='Price = ' + str(count.price) + ' $')
163     window.mainloop()
164
165 # Function button Yes

```



```

159 def button_yes():
160     # Select Score column and update to 1
161     my_df.data_f.iloc[count_object.currentObject,
162                        -2] = 1
163     # Condition for 10 products
164     if (count_object.currentObject + 1) % 10 == 0 :
165         # Compute classification score and update
166         dataframe
167         ml_score(my_df.data_f, my_df.data_f2)
168         update_df3(my_df.data_f, my_df.data_f2, var,
169                   var2, var3)
170
171     print(my_df.data_f)
172     # Update tkinter window
173     updateTk(var, var2, var3, count_object)
174
175 # Function button No
176 def button_no():
177     # Condition for 10 products
178     if (count_object.currentObject + 1) % 10 == 0:
179         ml_score(my_df.data_f, my_df.data_f2)
180         update_df3(my_df.data_f, my_df.data_f2, var,
181                   var2, var3)
182
183     print(my_df.data_f)
184     # Update tkinter window
185     updateTk(var, var2, var3, count_object)
186
187 # Function button stop -> close window
188 def button_stop():
189     window.destroy()
190
191 # Function button buy -> open chrome tav to the
192 # product URL
193 def button_buy():
194     driver = webdriver.Chrome(executable_path=
195                               any_var)
196     driver.get(my_df.data_f.iloc[count_object.
197                                currentObject, -4])
198
199 print(df2)
200 print(df3)
201 print(df2.size)
202 print(df3.size)
203
204 # 1st initialization tkinter window
205 window = tkinter.Tk()
206 # Get product picture from URL
207 im = Image.open(requests.get(count_object.img,
208                               stream=True).raw)
209 im = im.resize((400, 400), Image.ANTIALIAS)
210 # Create a PIL object and add to the canvas
211 photo = PIL.ImageTk.PhotoImage(image=im)
212 canvas = tkinter.Canvas(window, width=400, height
213                           =400)
214 canvas.pack()
215 canvas.create_image(0, 0, image=photo, anchor=
216                    tkinter.NW)
217
218 # Create labels
219 name = tkinter.Label(window, text =count_object.
220                       descr)
221 name.pack(anchor=tkinter.CENTER, expand=True)
222 price = tkinter.Label(window, text ='Price = ' + str
223                        (count_object.price) + ' $')
224 price.pack(anchor=tkinter.CENTER, expand=True)
225
226 # Create buttons
227 btn_yes = tkinter.Button(window, text="Yes", width
228                           =50, command=button_yes)
229 btn_yes.pack(anchor=tkinter.CENTER, expand=True)
230 btn_no = tkinter.Button(window, text="No", width=50,
231                          command=button_no)

```

```

219 btn_no.pack(anchor=tkinter.CENTER, expand=True)
220 btn_stop = tkinter.Button(window, text="Stop", width
221                            =50, command=button_stop)
222 btn_stop.pack(anchor=tkinter.CENTER, expand=True)
223 btn_buy = tkinter.Button(window, text="Buy", width
224                           =50, command=button_buy)
225 btn_buy.pack(anchor=tkinter.CENTER, expand=True)
226 window.mainloop()

```

## REFERENCES

- [1] <https://www.statista.com/statistics/991947/amazons-annual-technology-and-content-expenses/> , last consulted the 22 may 2021
- [2] Schmitt, Josephine & Debbelt, Christina Schneider, Frank. (2017). Too much Information? Predictors of Information Overload in the Context of Online-News Exposure. Information Communication and Society. 21. 1151–1167. 10.1080/1369118X.2017.1305427.
- [3] <https://amzscout.net/product-database/> , last consulted the 22 may 2021