# Using R to automate the classification of e-commerce products

Aidan Boland

4 November 2017

# Who are Clavis Insight?

# CLAVIS INSIGHT

- ▶ Global leader in action-ready insights for product manufacturers.

- ▶ We monitor online retailers to provide insights for manufacturers on their online performance.

Video link

# Brands



**Grocery, Personal, Home, & Pet**

**Health & Wellness**

**Beauty, Luxury, Fragrance**

**Toys & Games**

**Hardlines**

**Data & Standards**

**Alcoholic Beverages**

**Consumer Elec & Home Appliances**

**B2B**

**Sporting & Apparel**

CLAVIS INSIGHT

EARL

# Retailers

## 450 retailers across x countries

# Volume of Data Handled

- Clavis processes a large amount of data monthly.
    - $\approx$ 530,000,000 rows.
    - From 450 stores.
- $\approx$ 40,000 new products added to the database monthly.
- One particular bottleneck is classifying new products, adding tags to help with further analysis.

# Classification of Products

- eCommerce sites are fast moving and no two are exactly alike.
  - Different stores have different category names and structures.
- There is a need to provide a standardized category view across multiple sites.
  - Each manufacturer has a different view of categories.
  - Need to be flexible when assigning a product to a category.

# Classification of Products



- Amazon
  - Grocery & Gourmet Food > Candy & Chocolate > Bars
- Walmart
  - Food > Candy & Gum > Chocolate

# Classification of Products

| | Portfolio Availability | Price and Promotions | Ratings | Reviews | Share of Search | Menu | Total Content | Core Content | Extended Content |
|---|---|---|---|---|---|---|---|---|---|
| **All** | 44 % | 90 % | 61 % | 47 % | 36 % | 22 % | 54 % | 61 % | 92 % |
| **Breath Fresheners** | 55 % | 100 % | 72 % | 60 % | 16 % | 0 % | 31 % | 49 % | 88 % |
| **Candy - Chocolate** | 39 % | 93 % | 62 % | 45 % | 33 % | 23 % | 66 % | 73 % | 93 % |
| **Candy - Non-Chocolate** | 59 % | 83 % | 59 % | 41 % | 17 % | 20 % | 32 % | 41 % | 88 % |
| **Grocery** | 48 % | 88 % | 69 % | 50 % | 20 % | 41 % | 68 % | 75 % | 94 % |

CLAVIS INSIGHT

EARL
ENTERPRISE APPLICATIONS OF THE R LANGUAGE

# Classification of Products

- When a new product is found on a retailer, the information is collected.
- Each new product is then classified into it's correct category.
- As a manual task this is very time consuming, error prone and requires a lot of manual intervention by users.
  - $\approx$ 5 hours per customer end to end.
  - Source of complaints from customers.

- The difficulty of this process was a major indication for the need to automate.

CLAVIS INSIGHT

EARL
ENTERPRISE APPLICATIONS OF THE R LANGUAGE

# Automating Classification

# Roadmap for Automating Classification

- Research
  - Find suitable algorithm.
- Beta testing
  - Test the algorithm within the current process.
- Integrating code
  - Build the code into the current company platform.

EARL

- Task
  - Classify new products into their correct categories and subcategories (up to 8 levels).
    - Efficient.
    - Accurate.

- Solution
  - Supervised classification.

# Supervised Classification.

- Learn from the current products in order to classify new products.
- Many methods were considered (maxent, random forests, k-nn).
- Final algorithm was personalized to this particular task.

  - Classify's all levels of hierarchy simultaneously.
  - Relatively efficient ($<$30 seconds for 1000 products)
  - Ability to handle all languages.

- Researching and creating the algorithm is only one step.

# Beta Testing

## Shiny Application



**CLAVIS** INSIGHT

Dynamic Product Classification

If you are finding the tool slow, see instructions on how to run the tool on you local machine here.

The cpc input file should have the following two columns: 'trusted_rpc' and 'trusted_product_description' , and at least one of the following columns; 'manufacturer', 'brand', 'category', 'dimension1', 'dimension2', 'dimension3', 'dimension4', 'dimension5', 'dimension6', 'dimension7', 'dimension8'.

The harvest (candidate) file should have the following two columns: 'rpc' and 'Product_Description'. If 'Harvest_URL' and/or 'Product Image' are provided they will be displayed.

The probabilities are colour coded; red (too unsure to classify) , orange (low chance of being correct) , yellow (high chance of being correct) and green (likely to be correct) .

Report errors to Aidan Boland.

| Data | Hierarchical Structure |
|---|---|
| Manual Classification | Advanced |

**Load in cpc/training file**

| Browse... | Portfolio Availability-2016-11-15.csv |
|---|---|
| | Upload complete |

**Load in harvest/candidate file**

| Browse... | No file selected |
|---|---|

Perform Cross-validation

| Data Overview | Product Classification |
|---|---|

## Cross-validation

Lower percentages are better.

| manufacturer | brand | category | dimension1 | dimension2 | dimension3 | dimension4 | dimension5 | dimension6 | dimensio |
|---|---|---|---|---|---|---|---|---|---|
| 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% |

## Manufacturer

| Philips | 680 |
|---|---|

## Brand

| Philips Norelco | 286 |
|---|---|
| Philips Sonicare | 394 |

## Categories

CLAVIS INSIGHT

EARL
ENTERPRISE APPLICATIONS OF THE R LANGUAGE

# Beta Testing

- Shiny application
    - Hosted on an internal server.
- Allows the ability to manually load in training data and new data.
- Final downloaded is in the correct format to load into the production database.

- Still manual touch.

# Full Integration

# Integration into Platform

- The backend is coded in Java.
- Algorithm is coded in R.

- Needed to integrate the R code into the backend code.
- Must follow the companies standard practices.
  - Version control (bitbucket).
  - Change testing (unit tests).

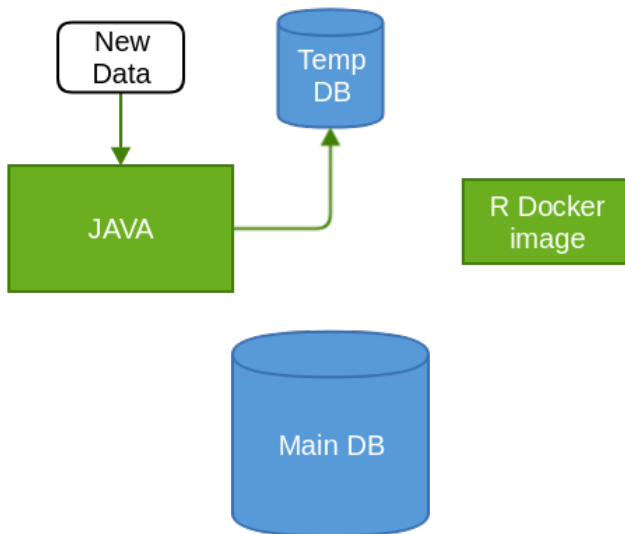CLAVIS INSIGHT

EARL

# Solution

- R library (including unit tests).
- Stored in a repository on Bitcucket.
  - This allows version control and updates which are in line with company practices.


- The **plumber** library is used to allow the code to be called using an API.
  - Arguments in the API point the R code to the database containing the new data.
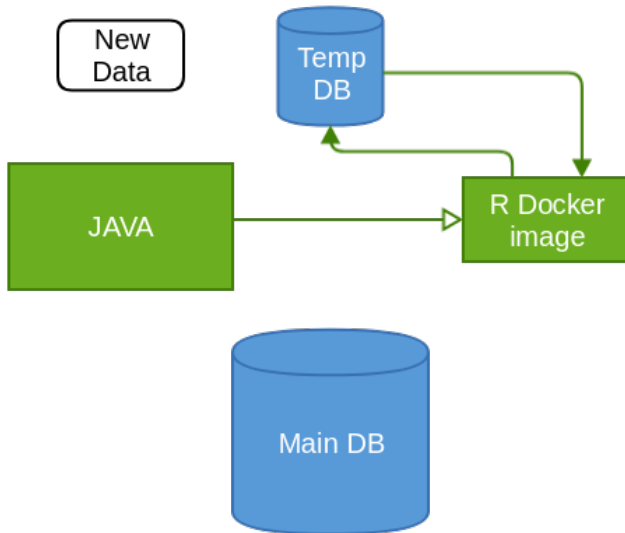  - R code classify's the new data and updates the database with the new categories.

EARL

# Solution

- A Docker image containing the R library and any requirements is stored.
  - The Docker image is automatically rebuilt when code is updated or a new production branch is created.
  - R library is recompiled and unit tests are run when the Docker image is rebuilt.

- Java code can spin up a new server and run the Docker image which will contain the code and will allow calls to be made to the API.
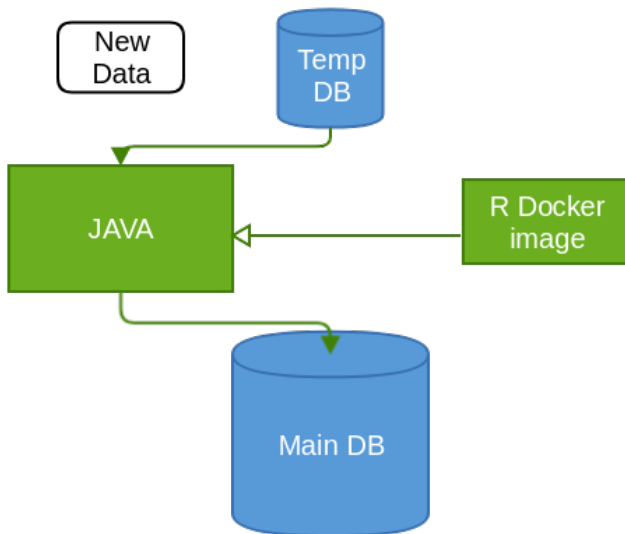  - Multiple Docker images can run at any one time.

EARL

# Final Flow of Data

# Final Flow of Data

# Final Flow of Data

# Improvement in Time to Add New Products

- Before automated classification
  - $\approx$ 5 hours per customer end to end.
  - 76 customers processed monthly.
  - Major source of complaint from customers.

- After automated classification
  - $\approx$ 2 hours per customer end to end.
  - 159 customers processed monthly.
  - No 'red' complaints since March.

CLAVIS INSIGHT

EARL