# Coffee Shops in San Francisco

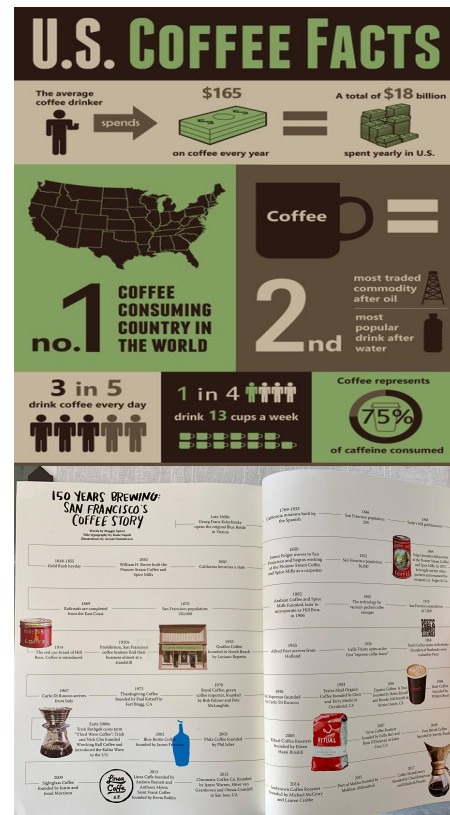**Looking for meaningful, unconventional and unique experiences**

# Introduction



- Coffee in the **U.S.**
  - The average coffee drinker consumes **3.1 cups of coffee daily**
  - There are about **100 million coffee drinkers** in the U.S.
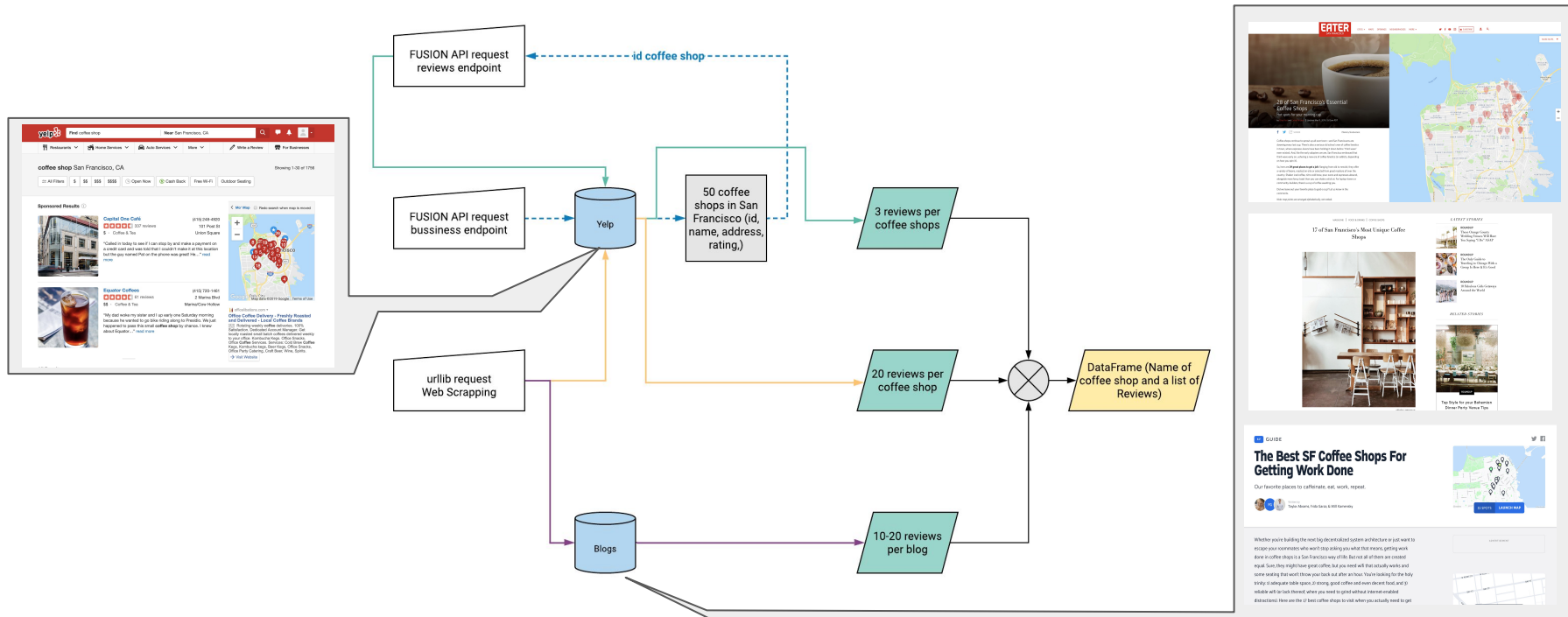  - American workers spend about **$20 per week** on coffee

- Why **San Francisco**?

  "()... you'll find an impressive number of coffee startups in San Francisco, and a new generation of tech-obsessed coffee-drinkers reshaping the city with every blink." (**Drift, Volume 7: San Francisco, July 15, 2018**)
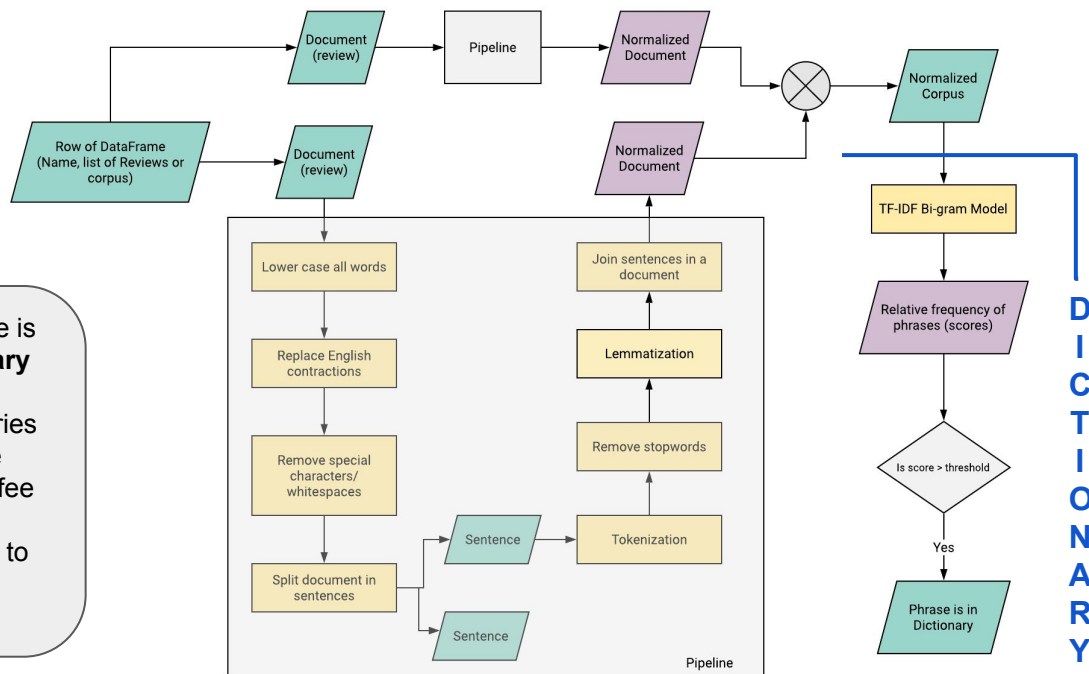
- Which elements define a pleasurable or disappointing experiences?
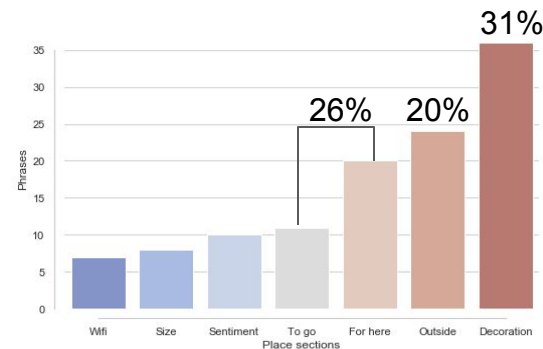
# Acquisition



FUSION API request reviews endpoint

FUSION API request bussiness endpoint

id coffee shop

Yelp

50 coffee shops in San Francisco (id, name, address, rating$_i$)

3 reviews per coffee shops

20 reviews per coffee shop

urllib request Web Scrapping

Blogs

10-20 reviews per blog

DataFrame (Name of coffee shop and a list of Reviews)
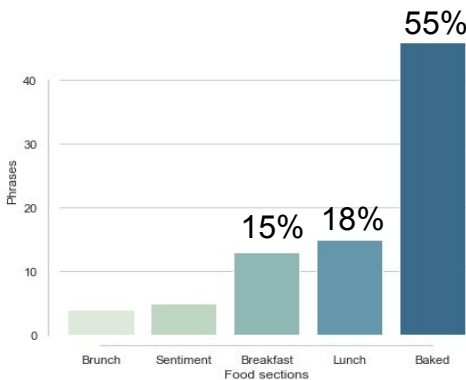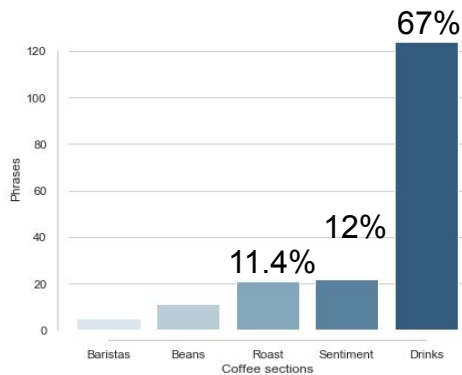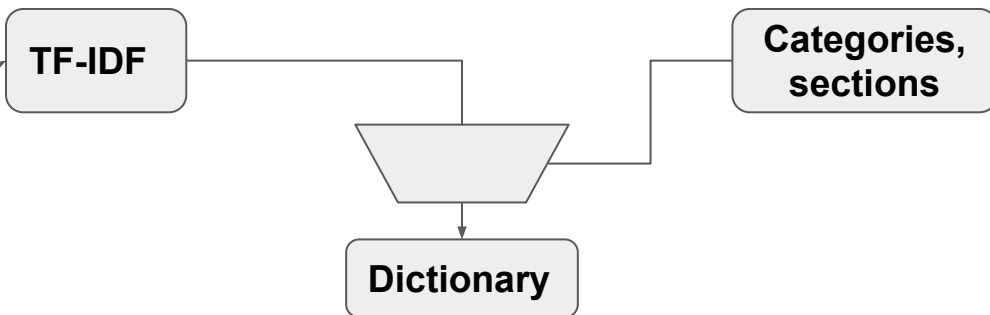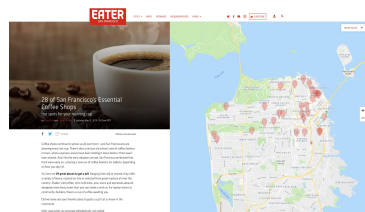
# Wrangling and Pre-processing



The preprocessing pipeline is used to to build a **Dictionary of key-phrases**. They are labeled in different categories and sections related to the experience around the coffee shop.
The same pipeline is used to work with **normalized reviews** later.

Why **Bi-grams**? We are looking for expressions as: **Pour over**, **latte art, cold brew**, **best coffee**, **amazing view**, **quiet place**, **seat available**, for instance.
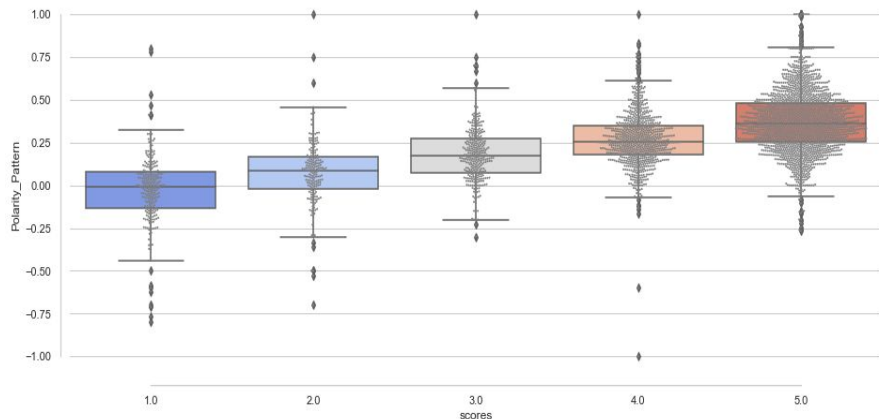
A **threshold** is used for avoiding rescue expressions with lower relative frequency.

Document (review)

Pipeline

Normalized Document

Normalized Corpus

Row of DataFrame (Name, list of Reviews or corpus)

Document (review)

Normalized Document

TF-IDF Bi-gram Model

Relative frequency of phrases (scores)

Is score > threshold

Yes

Phrase is in Dictionary

DICTIONARY

Pipeline

Lower case all words

Replace English contractions

Remove special characters/ whitespaces

Split document in sentences

Sentence

Sentence

Tokenization

Remove stopwords

Lemmatization

Join sentences in a document
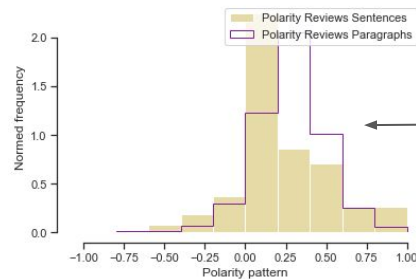
# Initial findings: Analysis of the Dictionary

# Initial findings: Sentiment Analysis of reviews

Could we measure how much **positive**, **negative** or **neutral** is the information of customer reviews?
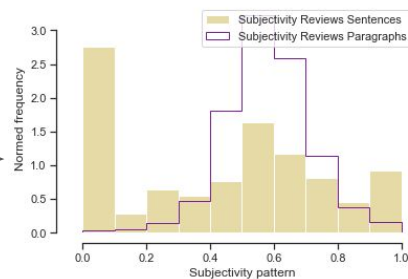
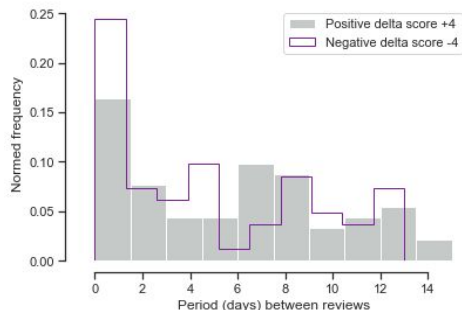How much the **polarity** and **subjectivity** change splitting reviews into sentences?



The **mean shifted at the left** when we split reviews.

**Change the shape** of the distribution: more neutral and subjective information.

# Initial findings: Period between reviews

How long take to write the next review for a specific coffee shop?
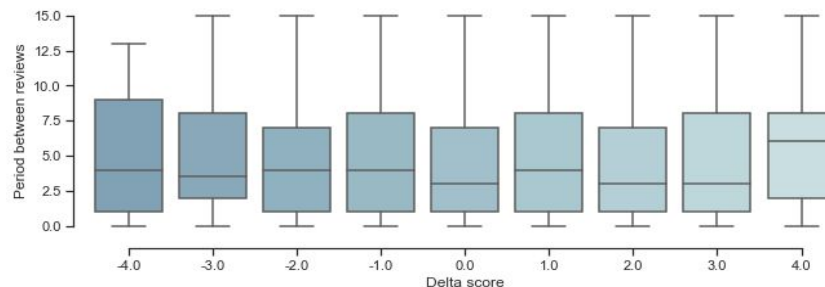


The difference between scores is **delta score** and the difference in days as the **period between reviews**

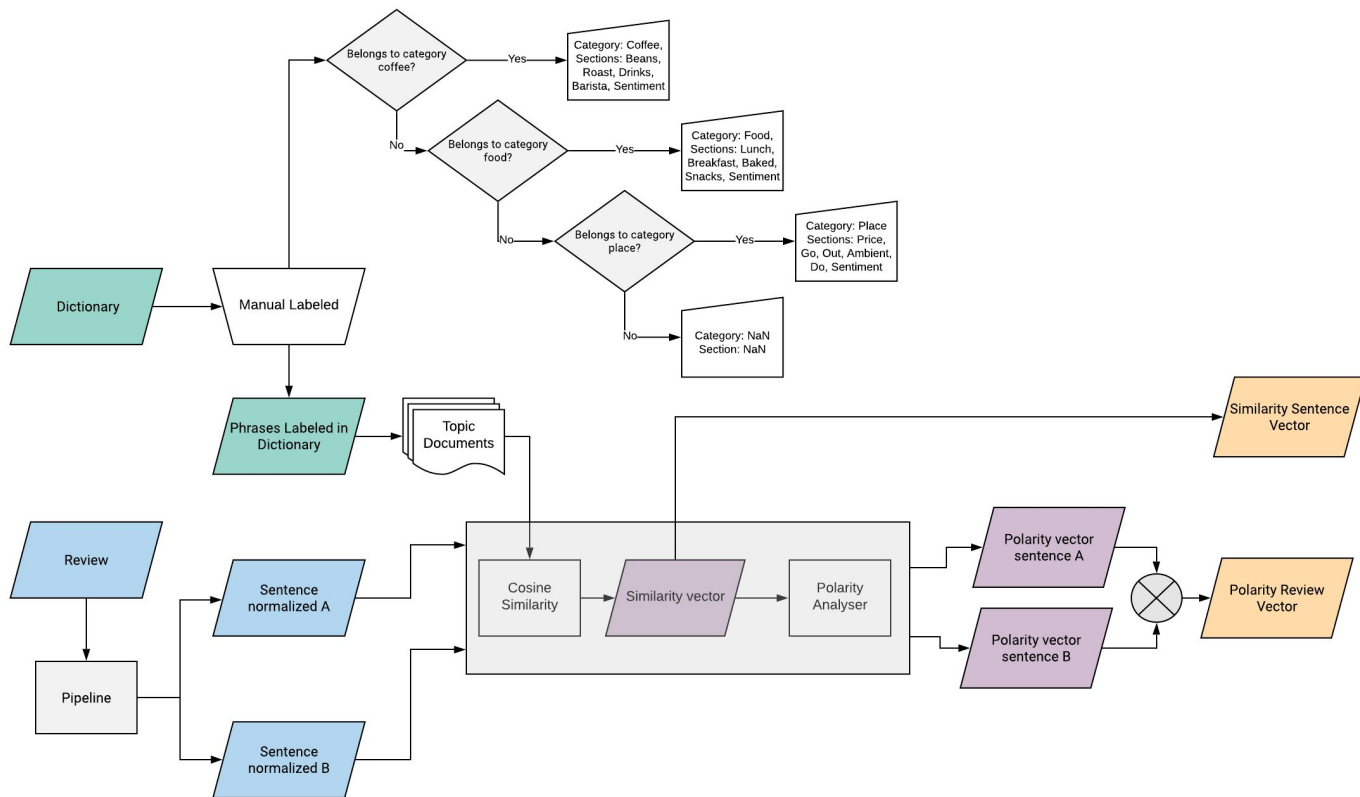The mean for **positive delta score is 5.6 days** and for **negative is 4.92 days**

Period between all possible Delta Score:

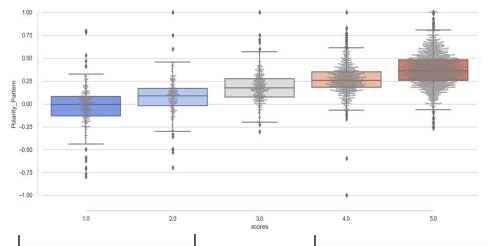**Positive Delta Score** of **+4** takes **more time** than the rest of delta scores

# SML: How feature vectors are built?

# SML: Selecting features and output vectors

**Output vectors**



**Negative reviews**
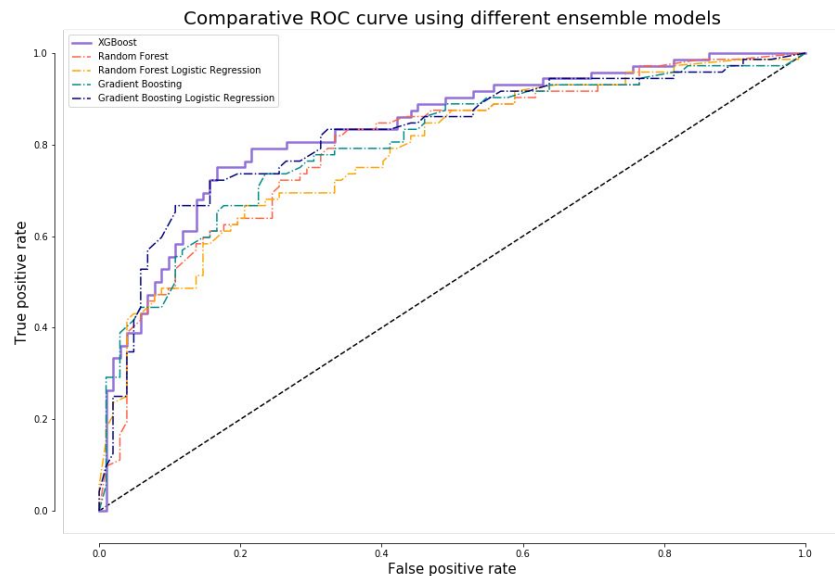
**Positive reviews**

**Selection of best feature vectors**

**6 features more strongly** related to the output variable

Statistical **test chi-squared** with a **confidence level of 5%** (SelectKBest)

| ID | SPECS | SCORE | ID | SPECS | SCORE |
|----|-------|-------|----|-------|-------|
| 15 | food sentiment | 5.413 | 9 | price | 1.068 |
| 10 | place sentiment | 3.445 | 3 | barista | 1.065 |
| 4 | coffee sentiment | 2.620 | 2 | drinks | 0.850 |
| 13 | breakfast | 2.499 | 5 | go | 0.815 |
| 8 | ambient | 2.215 | 14 | snacks | 0.663 |
| 11 | baked | 1.955 | 12 | lunch | 0.562 |
| 6 | do | 1.552 | 0 | beans | 0.165 |
| 7 | out | 1.171 | 1 | roast | 0.086 |

# SML: Training, tunning and testing models



Comparative ROC curve using different ensemble models

| Model | AUC | Accuracy |
|---|---|---|
| **XGBoost** | **83.46%** | **77.01%** |
| Random Forest | 79.71% | 71.84% |
| Logistic Regression Random Forest | 78.50% | 71.26% |
| Gradient Boosting | 80.22% | 75.29% |
| **Logistic Regression Gradient Boosting** | **82.16%** | **73.56%** |

| XGBoost | TP Reviews | TN Reviews |
|---|---|---|
| PP Review | 77 | 25 |
| NP Review | 15 | 57 |

# UML: Agglomerative Clustering using all features

**Bi-dimensional representation (PCA)**

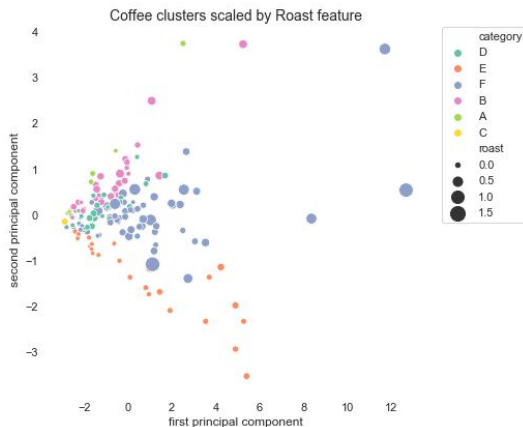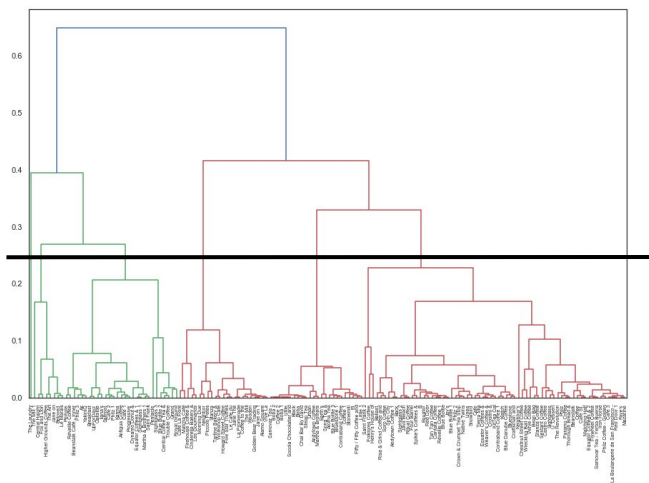How **features** are **determining** these clusters?



Coffee features

Food features

# UML: Dendrograms splitting features

At distance of **0.25** coffee features divide coffee shops in **6 groups**



Coffee clusters scaled by Roast feature

**Cluster F** higher roasting and beans sum, followed by B and D. Clusters A and E, is the lowest.

**Cluster A**, Le Marais, Art's Cafe, some bakeries among others



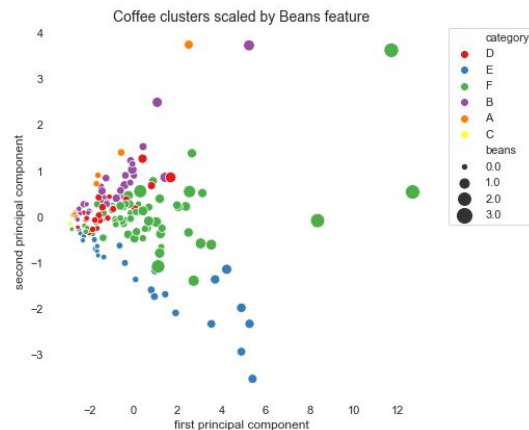Coffee clusters scaled by Beans feature

**Cluster F**: Andytown, Mazarine, Four Barrel, Red Door, Ritual Coffee, Saint Frank, Sightglass, Equator, Blue Bottle, Wrecking Ball, among others.

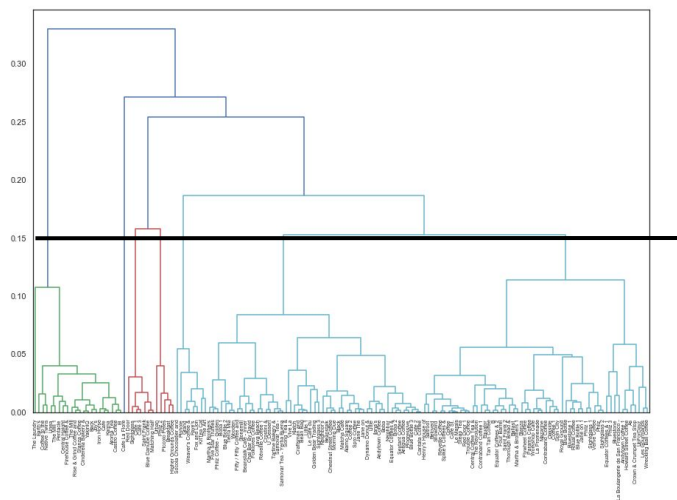**Cluster E**, Mostly tea and chai stores and chocolatier business

# UML: Dendrograms splitting features
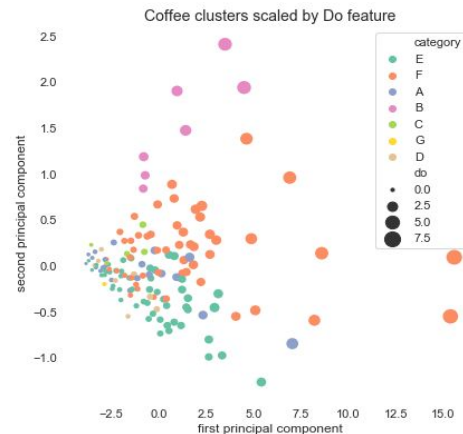
At distance of **0.15** coffee features divide coffee shops in **7 groups**

**Cluster B** higher Do values (feature about what to do in coffee shops): Matching Half, Red Door, Saint Frank, Sightglass (study, work)

**Cluster C** one of the lowest Do values. Piccolo Petes Urban, Higher Ground Coffee (bar styles, going to share with friends)



Coffee clusters scaled by Do feature

# Conclusion

- **Blogs put more attention in:**
  - Drinks, baked items and decoration in their reviews.

- **Features most strongly vinculated with the score of customer reviews:**
  - **Sentiments** associated with **drinks**, **place** and **food**, **breakfast**, **ambient**, **outside**, **baked elements**, **baristas**, **available of seats**, **what to do** and **price**.
  - Less relationship: **quality** and **variety of beans**, **roasting,** little **snacks** and **lunch** options.

# Conclusion

- **Supervised Learning:**
  - We evaluated different **Decision Trees** models to predict sentiment patterns in customer reviews. **XGBoost** got the best performance, with an AUC of 83.4% and an accuracy of 77 %.

- **Unsupervised Learning:**
  - Using **Agglomerative Clustering** and **dendrograms** as visualization tool, we found clusters with coffee shops distinguishing styles, what to do there and how much interested is the people in talking about **beans**, **roasting**, **drinks, places** and **food**.

# Resources

- Yelp API (FUSION) and web scraping from Yelp
- "*28 of San Francisco's Essential Coffee Shops: Hot spots for your morning cup*" (Ellen Fort and Caleb Pershan). Available here.
- "*17 of San Francisco's Most Unique Coffee Shops* (Katie Bush). Available here.
- "*The Best SF Coffee Shops For Getting Work Done*" (Taylor Abrams, Frida Garza, and Will Kamensky). Available here.
- Coffee gives me superpowers (Ryoko Iwata). Published on April 7, 2015
- DRIFT San Francisco (A. Goldberg, Velasco, Lee, E. Goldberg and Spicer). Published on July 15, 2018
- XGBoost: The Excalibur for Everyone (Raghu Raj Rai). Towards Data Science. Available here
- A Beginner's guide to XGBoost (George Seif). Towards Data Science. Available here
- Traditional Methods for Text Data (Dipanjan Sarkar). Towards Data Science. Available here
- Practical Statistics for Data Scientist (Peter Bruce and Andrew Bruce). O'REILLY, 2017.
- Text Classification is Your New Secret Weapon (Adam Geitgey, Medium). Available here
- A Practitioner's Guide to Natural Language Processing (Part I) - Processing and Understanding Text (Dipanjan Sarkar). Towards Data Science. Available here