# Capstone project

*Hitting the top notes*

Modeling on fragrance notes to classify ratings

**Allison Bishop**

# Problem statement

Are fragrance notes good predictors of average customer ratings?

# Background - My motivation

## CHERRIES—IN GENERAL

**Season:** late spring–late summer

**Taste:** sweet

**Weight:** light–medium

**Volume:** moderate

**Techniques:** flambé, poach, raw, stew

**Flavor Affinities**

cherries + almonds + cream + kirsch + vanilla

cherries + chocolate + walnuts

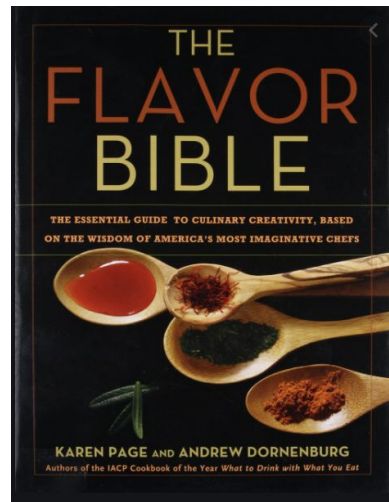cherries + coconut + custard

cherries + coffee + cream

cherries + goat cheese + ice wine vinegar + black pepper + thyme

cherries + honey + pistachios + yogurt

cherries + mint + vanilla

cherries + orange + sugar + dry red wine

cherries + sweet vermouth + vanilla



THE FLAVOR BIBLE

THE ESSENTIAL GUIDE TO CULINARY CREATIVITY, BASED ON THE WISDOM OF AMERICA'S MOST IMAGINATIVE CHEFS

KAREN PAGE AND ANDREW DORNENBURG

Authors of the IACP Cookbook of the Year *What to Drink with What You Eat*

# Background - Classes of notes

**Top**: Form initial impression. Selling point.
High volatility.
*Light, bright (like citrus fruits)*

**Middle**: Forms the body. 40-80% of total aroma.
Midrange volatility.
*Complex, midweight (like florals)*

**Base**: Foundation of fragrance. Brings depth.
Low volatility.
*Deep, heavyweight (like sandalwood)*

# Preprocessing - Dummify notes

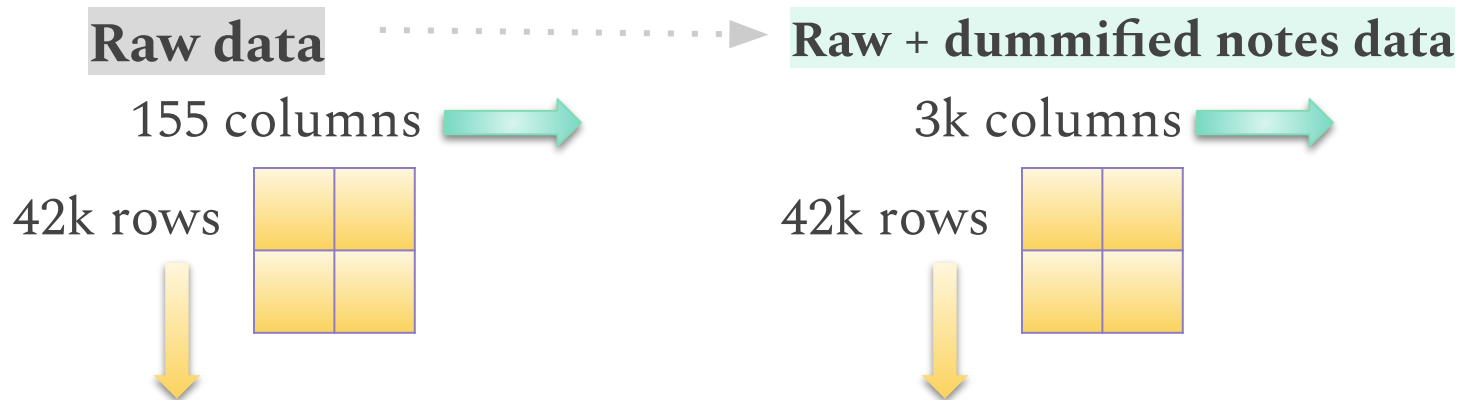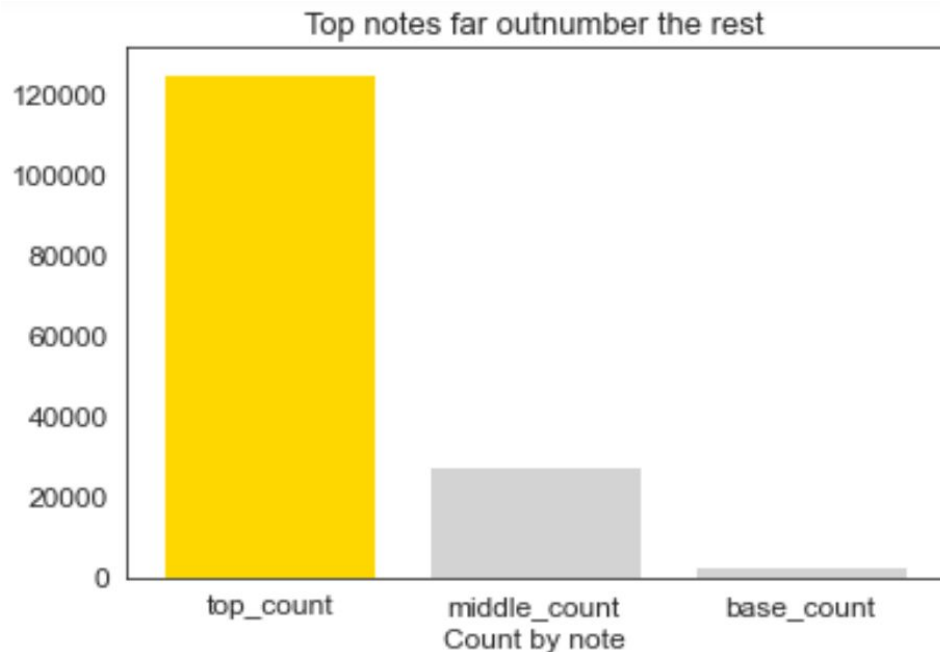| | title | 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|
| 0 | Aamal The Spirit of Dubai for women and men | Top0Turkish Rose | Top1Bulgarian Rose | Top2Bergamot | Top3Fruits |
| 1 | Aatifa Ajmal for women and men | Top0Nutmeg | Top1Rose | Top2Cumin | Middle0Amber |

*After*

| top_0_mandarin_orange | top_1_green_apple | top_2_thyme | middle_0_2_lavender |
|---|---|---|---|
| 1 | 1 | 1 | 1 |

# Preprocessing - Results

The dataset drastically changed shape

**Raw data** ┄┄┄┄┄▶ **Raw + dummified notes data**

155 columns

42k rows

3k columns

42k rows

# Data profile - Top, middle, base notes



Top notes far outnumber the rest

Abundance of top notes data

Project focus

*Quick hit (top) only*

# Data profile - Average ratings



Ratings approach normal distribution (some left skew)

After binning, modeled just on ratings 3 and 4

# EDA - Fragrance notes



Top 5 most frequent fragrance note categories
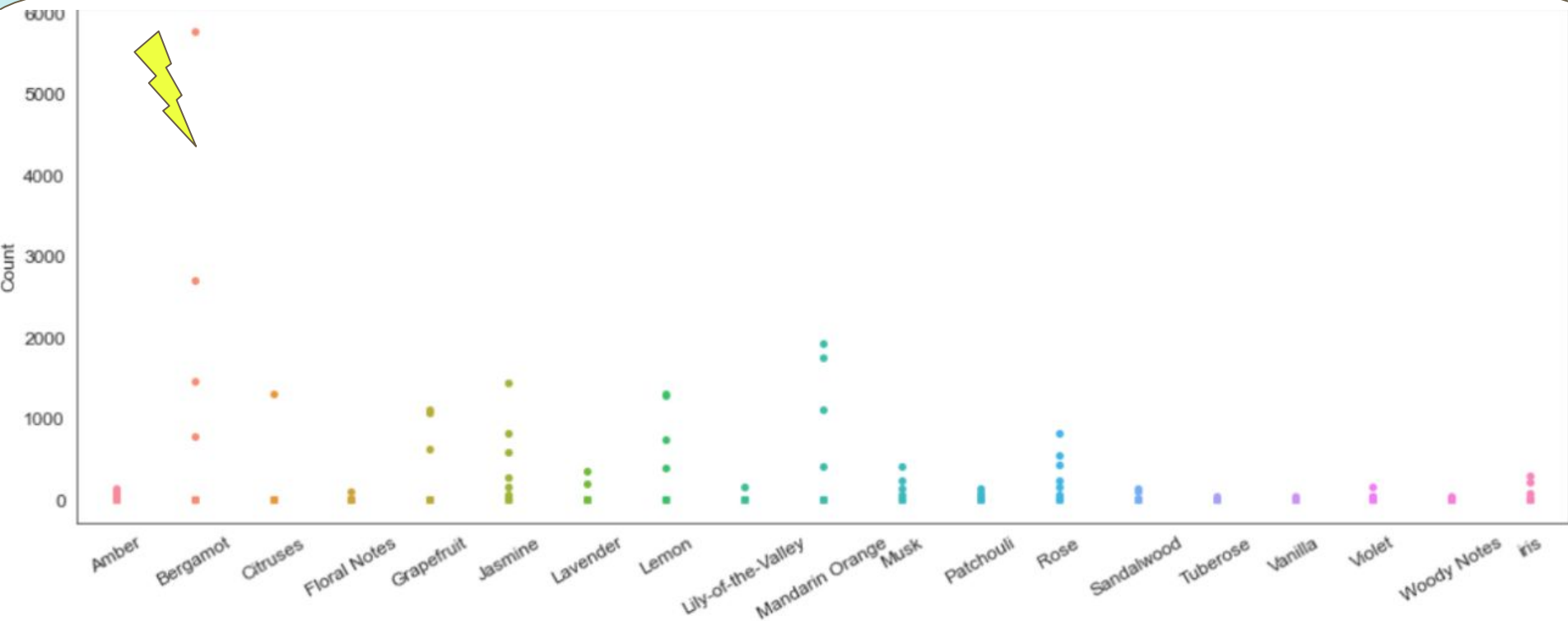
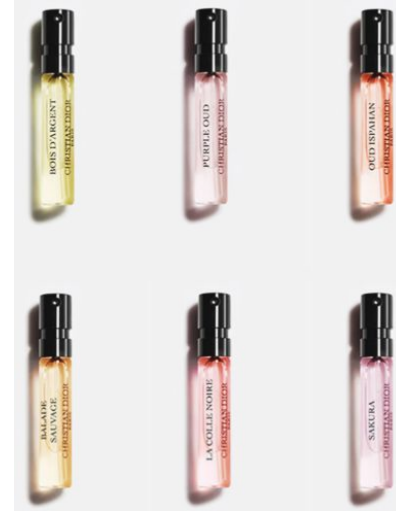# EDA - Counts of most frequently used notes

# **Sampling of data science methods**
*A bit like sampling fragrances!*

Logistic regression

Clustering

# Modeling results - Accuracy scores

Logistic regression
with _and_ without
Principal Component Analysis (PCA)

Baseline = 54%

Majority class = Rating 4

With PCA

Train set: 59.6%

Test set: 56.9%

40 features

Without PCA

Train set: 91.4%

Test set: 52.3%

# Modeling results - Accuracy scores

Tree-based models

|  | **Decision tree** | **Random Forest + GridSearchCV** |
|---|---|---|
| Baseline model = 54% | Train set: 92% | Train set: 92% |
| Majority class = Rating 4 | Test set: 54% | Test set: 60% |

# Modeling results: What we can infer



Strongest top notes indicators for ratings (logistic regression)
For ratings 3 and 4 on a scale of 5

# Clustering - Results (on fives ratings only)

**Silhouette score**

Cohesion *(intra-cluster distance)* - separation *(inter-cluster distance)*
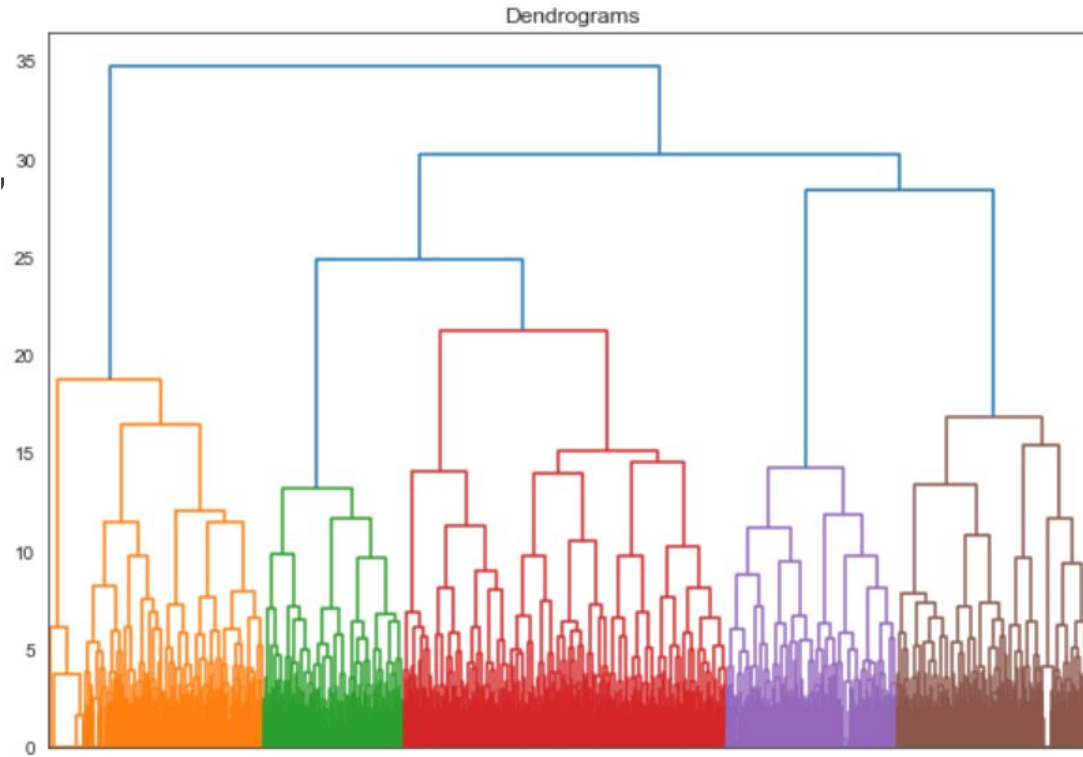
Range: -1 (worst) to 1 (best)

DBSCAN

Silhouette score: -0.3

*K*Means

Silhouette score: -0.2

# Clustering - Next steps



Feature
agglomeration
+
Hierarchical
clustering

# Conclusions and next steps (seriously)

Raw notes data has low predictive value

Tree-based system for product development

Live data stream

Now, a quick demo!