# WINE NOT?

## USING NLP TO EXPLORE FLAVOR PROFILES IN WINE REVIEWS

GENEVIEVE MCGUIRE

# SOMM BACKGROUND...

Goal: Develop a tool useful for both sommeliers and casual wine drinkers

# THE WORKFLOW

## DATA

WineEnthusiast
Magazine reviews from
2017
130k data points

## CLEANING

Domain stop words
Lemmatization
Remove hapaxes

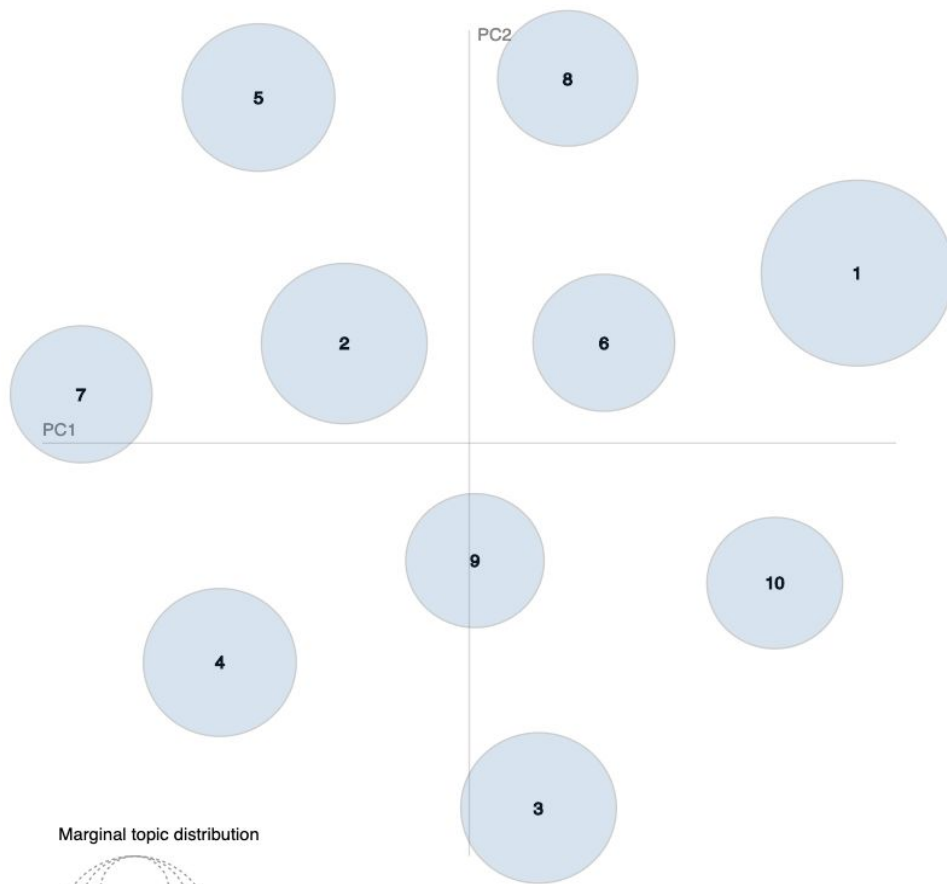## VECTORIZER

CountVectorizer
TF-IDF

## LSA & LDA

TruncatedSVD
NMF
LDA: Gensim
& Sklearn

kaggle™

# NMF TOPICS

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| **RED** | black cherry pepper spice licorice | | ripe rich structure wood balance | red fruity | soft cherry raspberry | dry tannic core herb | | blackberry dark chocolate blend tannic | body light medium color | berry plum herbal spice earthy |
| **ROSE** | | | | fruity bright fresh currant | sweet simple cherry raspberry | dry | | | strawberry light | |
| **WHITE** | | apple citrus white lemon crisp | | bright fresh | sweet | dry bone herb | oak vanilla toast barrel butter | | | |

## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribution

2%

5%

10%

## Top-30 Most Salient Terms[1]

| | 0 | 1,000 | 2,000 | 3,000 | 4,000 | 5,000 |

apple
peach
blackberry
medium
cherry
black
lemon
pear
crisp
lime
citrus
light
butter
pineapple
attractive
berry
herbal
meat
fruity
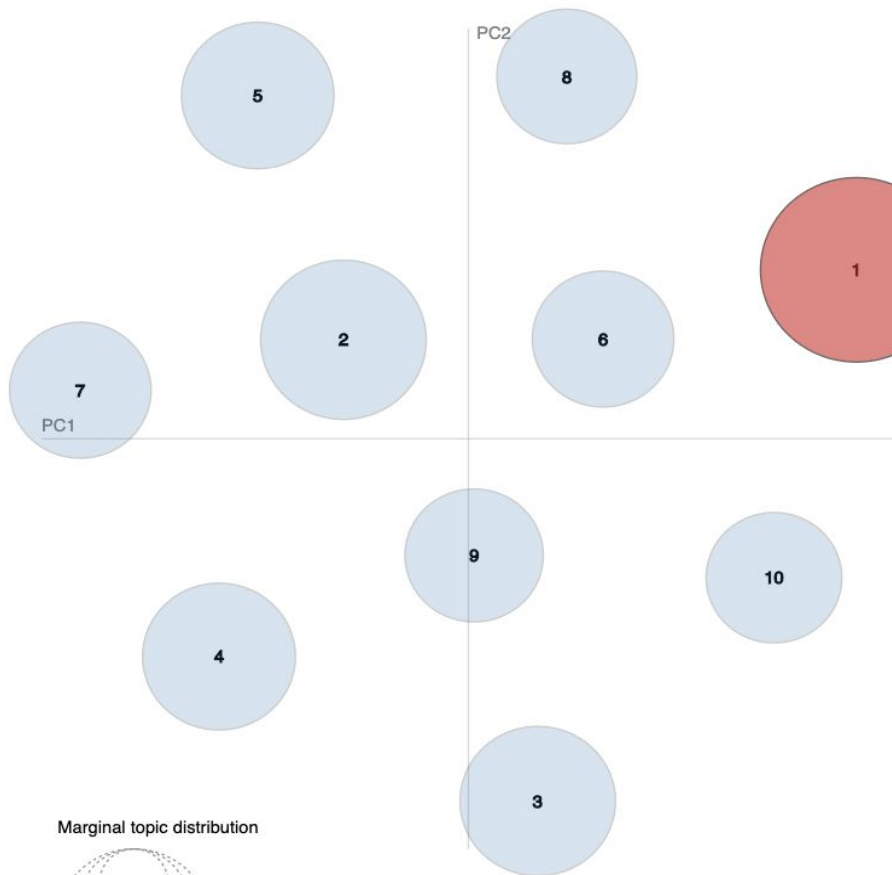strawberry
honey
oak
tropical
wood
barrel
pair
cranberry
old
raspberry
coffee

Overall term frequency

Estimated term frequency within the selected topic
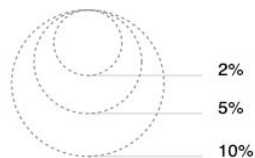
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)
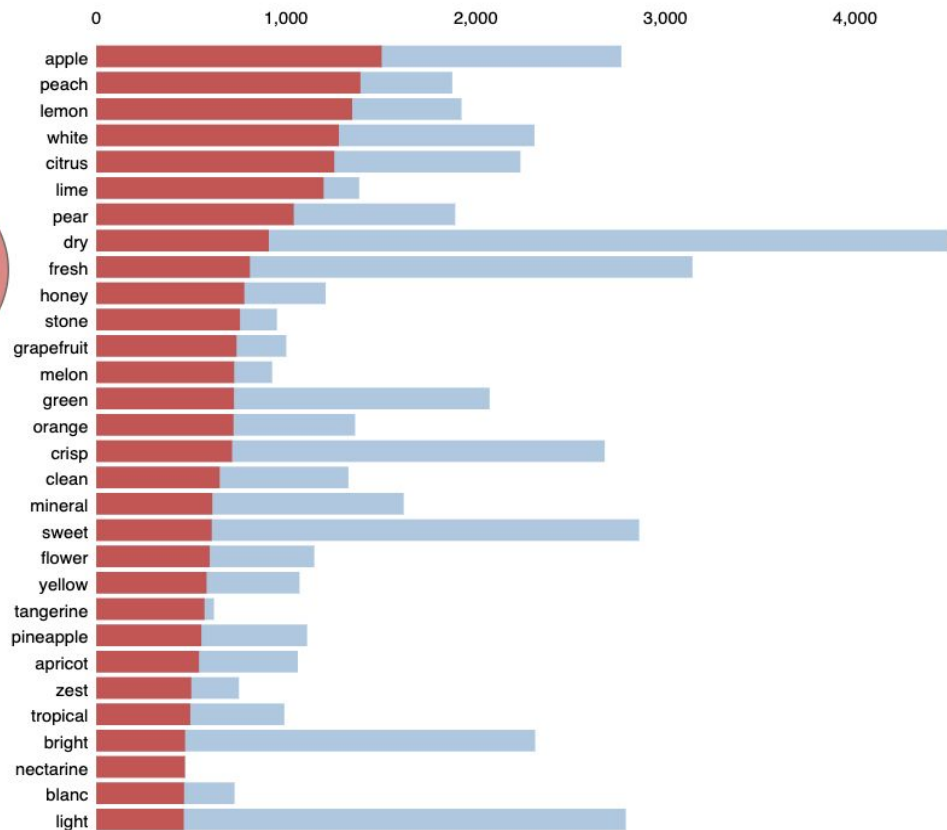
## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribution

2%

5%

10%

## Top-30 Most Relevant Terms for Topic 1 (15.9% of tokens)

| | 0 | 1,000 | 2,000 | 3,000 | 4,000 |

apple
peach
lemon
white
citrus
lime
pear
dry
fresh
honey
stone
grapefruit
melon
green
orange
crisp
clean
mineral
sweet
flower
yellow
tangerine
pineapple
apricot
zest
tropical
bright
nectarine
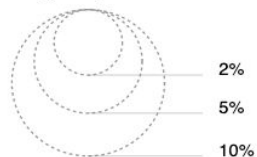blanc
light

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)
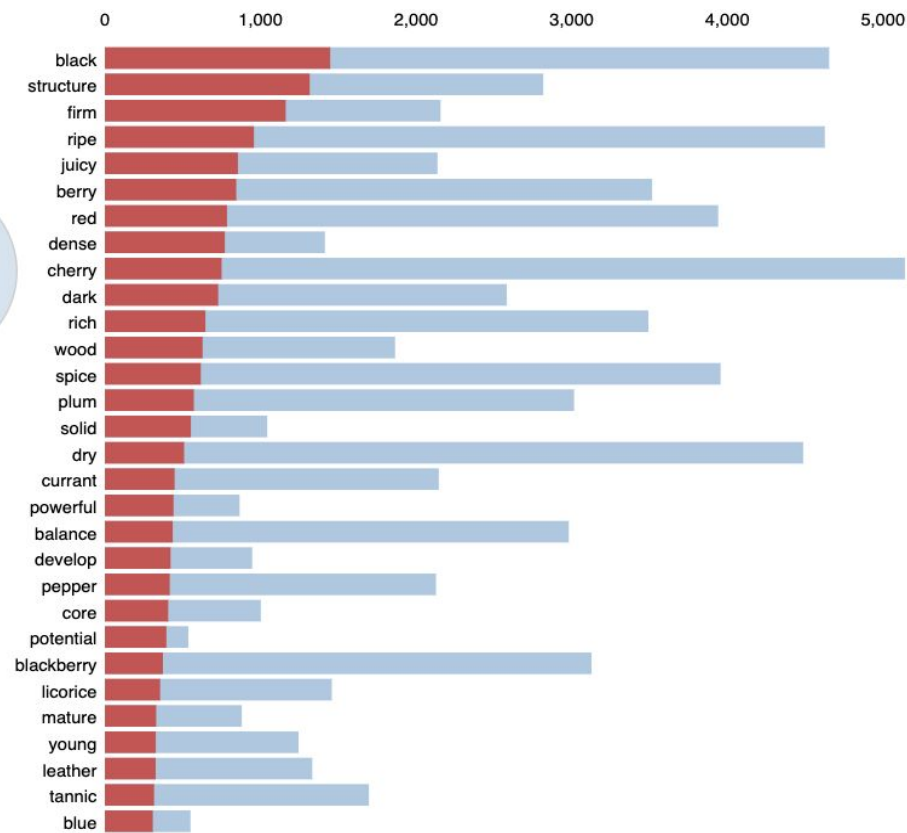
## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribution

2%

5%

10%

## Top-30 Most Relevant Terms for Topic 2 (11.8% of tokens)

black
structure
firm
ripe
juicy
berry
red
dense
cherry
dark
rich
wood
spice
plum
solid
dry
currant
powerful
balance
develop
pepper
core
potential
blackberry
licorice
mature
young
leather
tannic
blue

Overall term frequency

Estimated term frequency within the selected topic
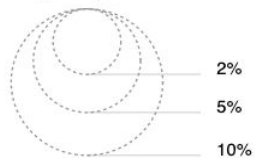
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

## Intertopic Distance Map (via multidimensional scaling)

## Top-30 Most Relevant Terms for Topic 3 (10.4% of tokens)

Marginal topic distribution

2%
5%
10%

Overall term frequency
Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

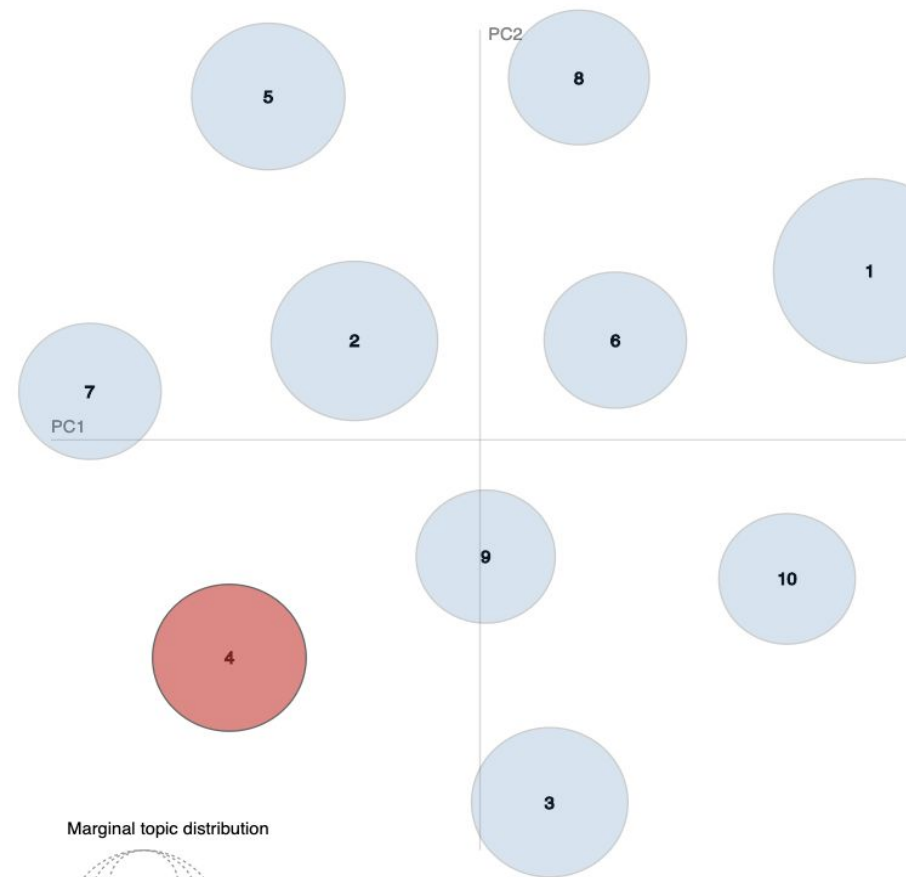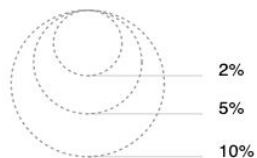## Intertopic Distance Map (via multidimensional scaling)



PC2

5

8

1

7

2

6

PC1

9

10

4

3

**Marginal topic distribution**

2%

5%

10%

## Top-30 Most Relevant Terms for Topic 4 (10.1% of tokens)



| | 0 | 1,000 | 2,000 | 3,000 | 4,000 | 5,000 |
|---|---|---|---|---|---|---|

black
cherry
dark
coffee
blackberry
chocolate
espresso
oak
plum
spice
licorice
dry
blueberry
vanilla
pepper
cassis
toast
cedar
blend
purple
deep
herb
boysenberry
sage
tobacco
body
lead
ripe
graphite
red

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)
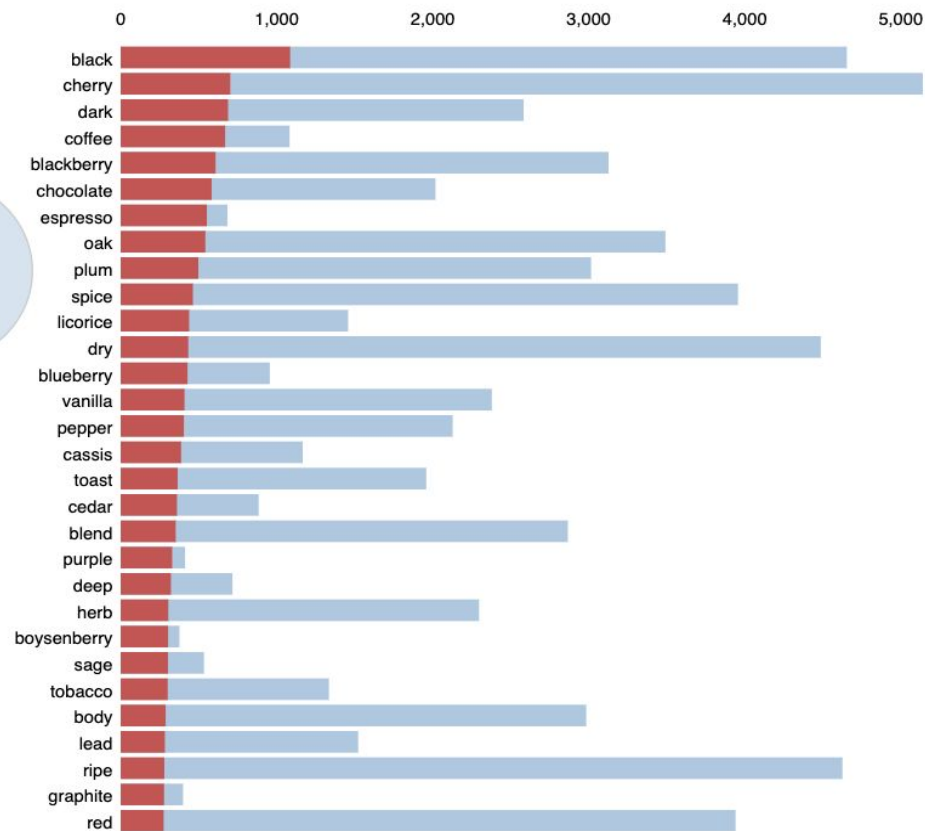
## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribution

2%

5%

10%

## Top-30 Most Relevant Terms for Topic 5 (10% of tokens)

Overall term frequency
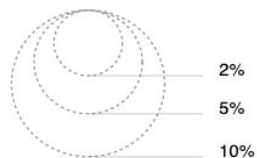
Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

## Intertopic Distance Map (via multidimensional scaling)
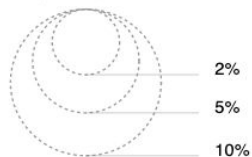
PC2

PC1

Marginal topic distribution

2%

5%

10%

## Top-30 Most Relevant Terms for Topic 6 (8.6% of tokens)

| | 0 | 1,000 | 2,000 | 3,000 | 4,000 | 5,000 |
|---|---|---|---|---|---|---|

toast
rich
oak
butter
wood
ripe
vanilla
barrel
old
pineapple
vine
richness
tropical
apple
creamy
pear
balance
caramel
spice
sweet
apricot
honey
opulent
ferment
french
toasty
peach
cream
body
oaky

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)
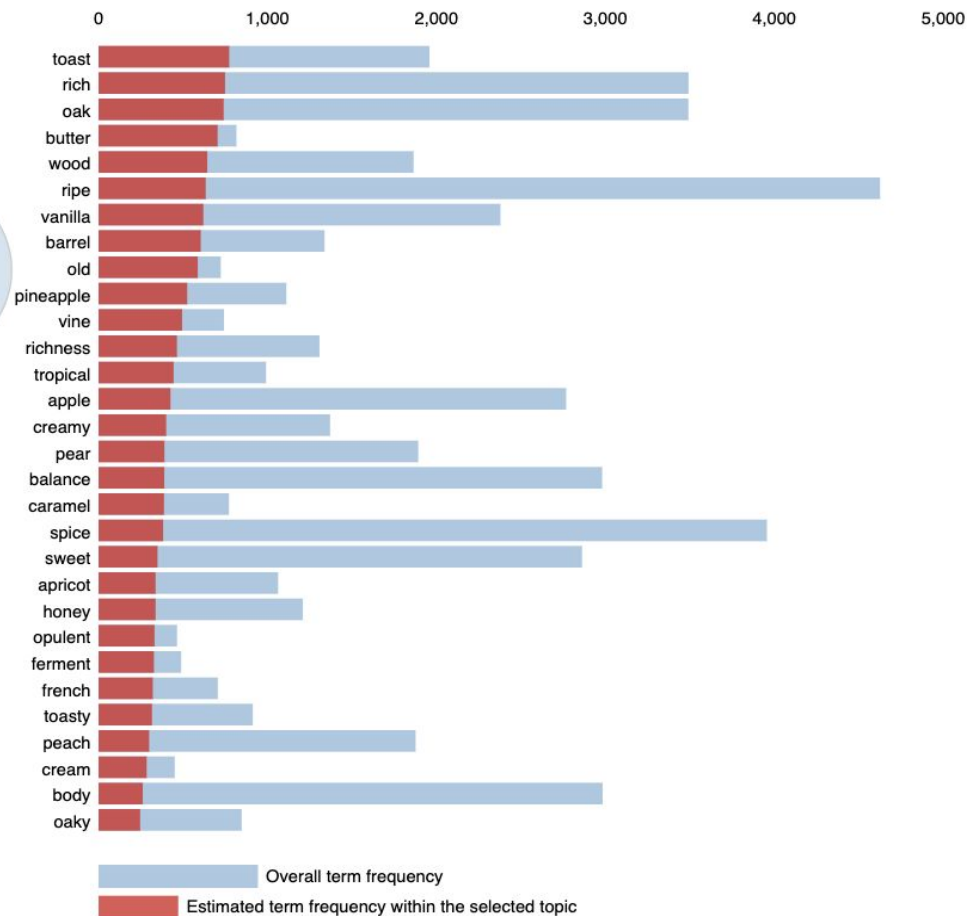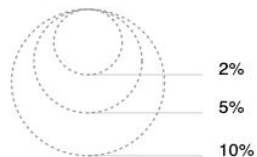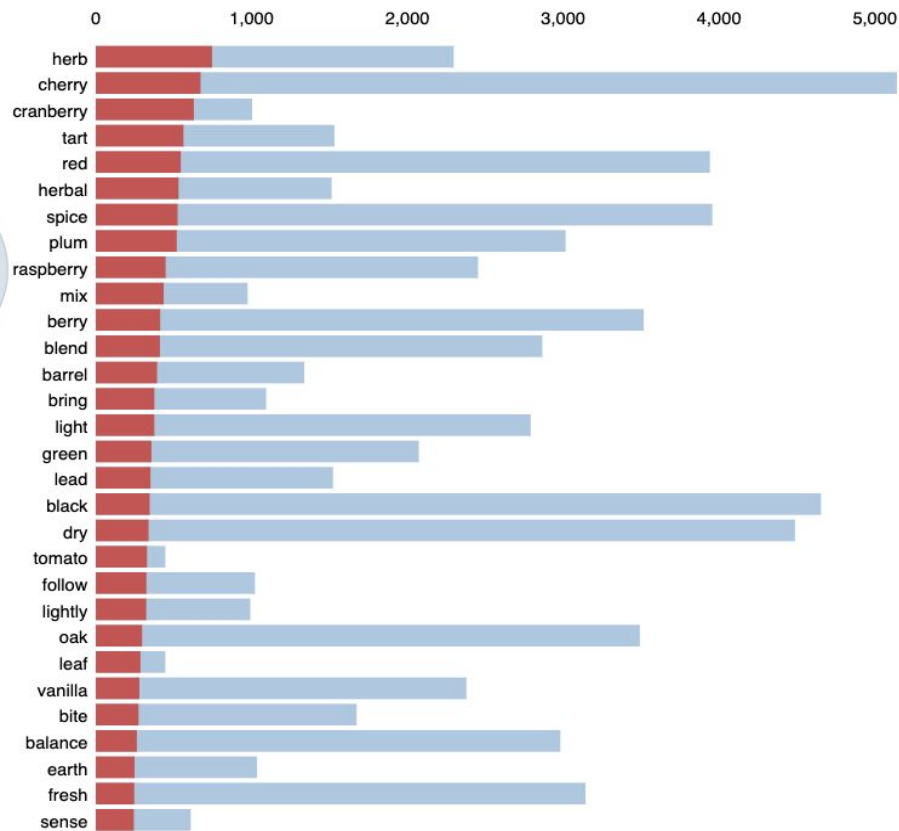
## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribution

2%

5%

10%

## Top-30 Most Relevant Terms for Topic 7 (8.6% of tokens)

| | 0 | 1,000 | 2,000 | 3,000 | 4,000 | 5,000 |
|---|---|---|---|---|---|---|

herb
cherry
cranberry
tart
red
herbal
spice
plum
raspberry
mix
berry
blend
barrel
bring
light
green
lead
black
dry
tomato
follow
lightly
oak
leaf
vanilla
bite
balance
earth
fresh
sense

Overall term frequency

Estimated term frequency within the selected topic
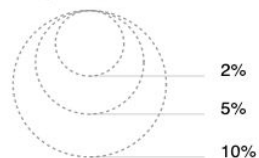
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)
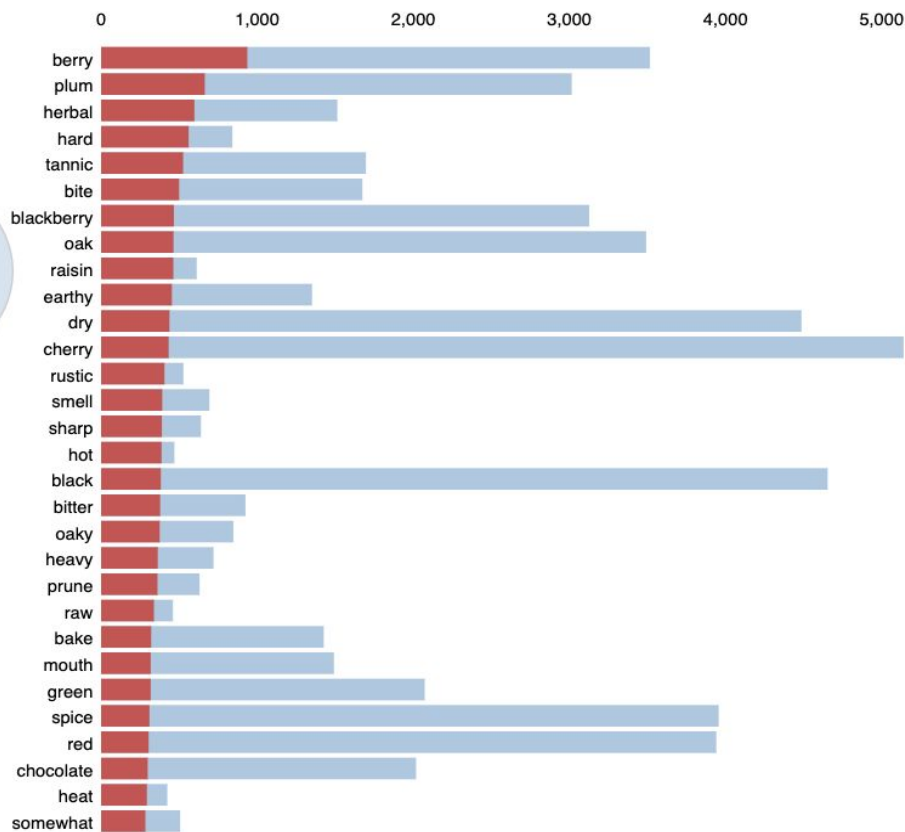
# Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribution

2%

5%

10%

# Top-30 Most Relevant Terms for Topic 8 (8.4% of tokens)

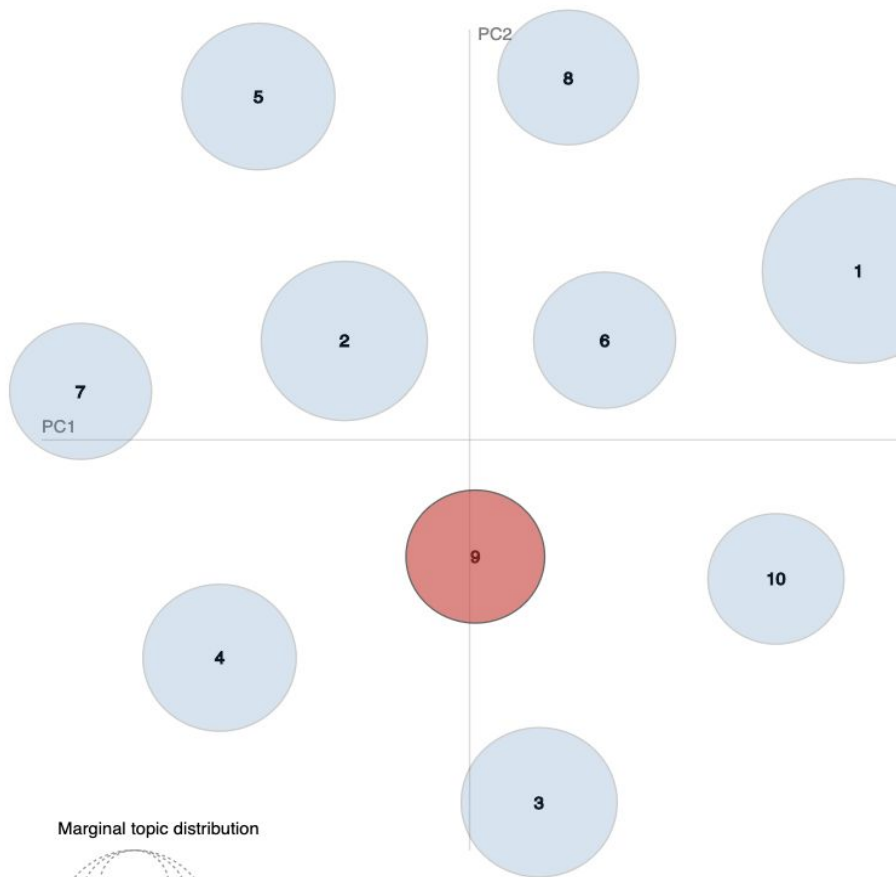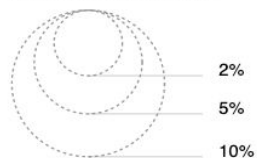| | 0 | 1,000 | 2,000 | 3,000 | 4,000 | 5,000 |
|---|---|---|---|---|---|---|
| berry | | | | | | |
| plum | | | | | | |
| herbal | | | | | | |
| hard | | | | | | |
| tannic | | | | | | |
| bite | | | | | | |
| blackberry | | | | | | |
| oak | | | | | | |
| raisin | | | | | | |
| earthy | | | | | | |
| dry | | | | | | |
| cherry | | | | | | |
| rustic | | | | | | |
| smell | | | | | | |
| sharp | | | | | | |
| hot | | | | | | |
| black | | | | | | |
| bitter | | | | | | |
| oaky | | | | | | |
| heavy | | | | | | |
| prune | | | | | | |
| raw | | | | | | |
| bake | | | | | | |
| mouth | | | | | | |
| green | | | | | | |
| spice | | | | | | |
| red | | | | | | |
| chocolate | | | | | | |
| heat | | | | | | |
| somewhat | | | | | | |

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
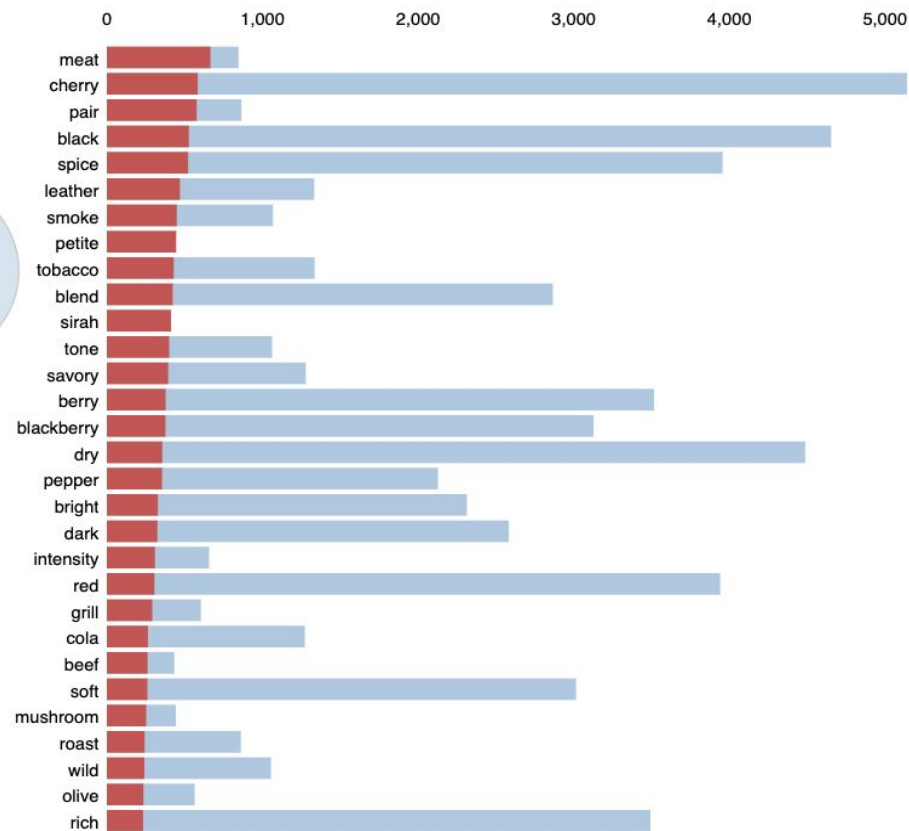2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

Intertopic Distance Map (via multidimensional scaling)

Top-30 Most Relevant Terms for Topic 9 (8.2% of tokens)

Marginal topic distribution

- 2%
- 5%
- 10%

Overall term frequency

Estimated term frequency within the selected topic
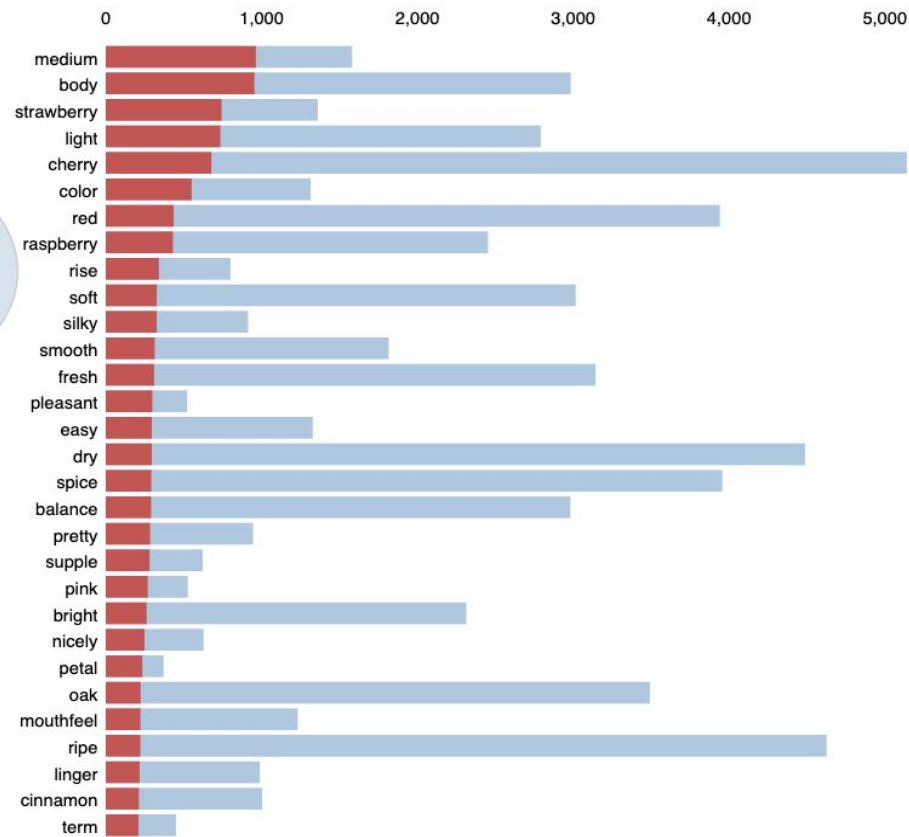
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

## Intertopic Distance Map (via multidimensional scaling)

## Top-30 Most Relevant Terms for Topic 10 (7.9% of tokens)

Marginal topic distribution

2%

5%

10%

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

# CONCLUSIONS...

Sommeliers' jobs are safe! (For now...)
Basic clusters can't beat a good somm's taste-memory, but could help a casual wine drinker find more of what they like

# FUTURE WORK

**01** REFINE CLUSTERS

Go deeper into crucial vocabulary

**02** RECOMMENDATION TOOL

Suggestions based on flavor groups

**03** INSTASOMM

Make use of labelled data in dataset

# THANK YOU!

## QUESTIONS?

gmmcguire2@gmail.com