



# An exploration of Wine Reviews and the Big Data of wines

Giacomo Freccero, DSI General Assembly 2021

# The Process:

- Some wine notions
- Goals
- Data Collection & EDA
- Findings
- Limits and potential improvements



# WINE 101

© WINE FOLLY



Grape  
Variety:  
Merlot



Grape  
Variety:  
Chardonnay



Wine

Wine

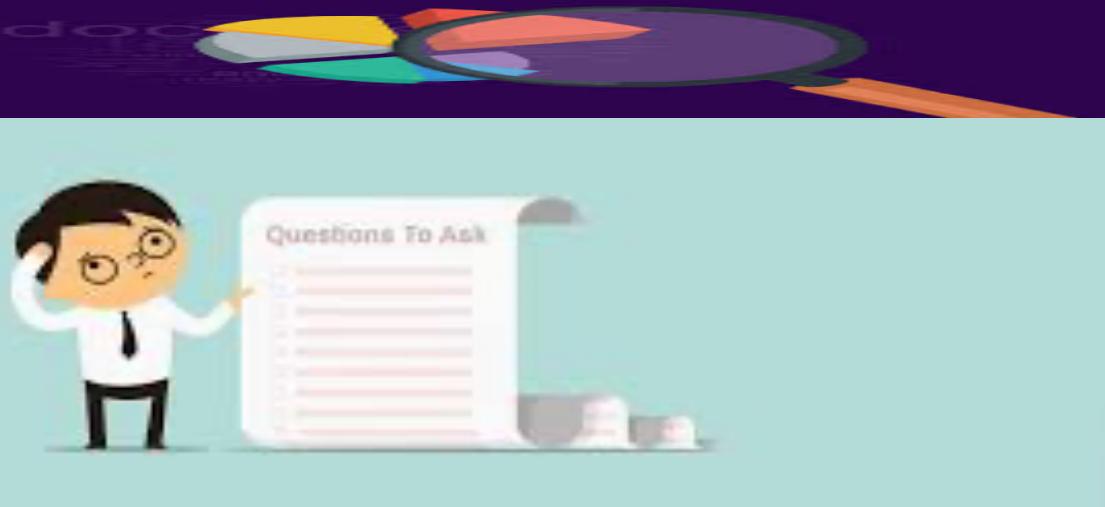


- B: Balance (fruit vs acidity)
- L: Length
- I: Intensity
- C: Complexity - 3 layers:
  - a) Aromas and Flavors of the grape and alcoh. fermentation
  - b) post-fermentations
  - c) maturation

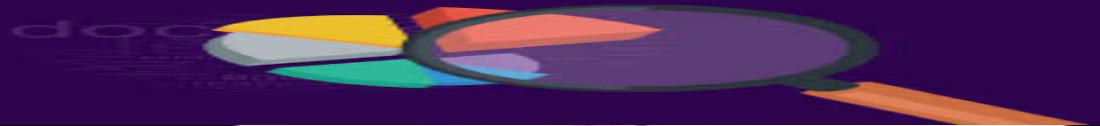


Some wine notions

# Goals



- Can we predict grape varieties from the text reviews & what words are the best predictors?
- How are the grape varieties connected in terms of common flavor characteristics? Can we create a map of flavors?
- Can we predict above or below median rating from the text reviews and what words are the best predictors?
- Can we predict ratings from text reviews?
- Can we predict value for money from the text reviews?
- Can we predict if the wine will be good, very good, superb or excellent from the text reviews?
- Can we predict Rating from the winery?
- Can we predict above/below median rating from the appellation?



## Why should I care?



- Create dictionaries of flavors that are non biased for each gv
- Get descriptors that are an indication of quality
- Understand what flavor characteristics gv share, mapping them and start a wine recommender. E.g. if you liked Californian Cab you will like Italian Aglianico.
- Which flavors are good value for money and which are expensive to buy
- Which flavors are associated with exceptional wines, which with bad wines?
- Which wineries consistently produce high rated wines (good investments)
- Which appellation is likely to predict above average wines
- Who would care? Wine Education/Wine Clubs/Wholesalers/Retail /Wine Investors/Wine consumers

# Data Collection

[41,894 total results]

RELATED REVIEWS 21-40 of 41,894

Filter by ▼ Sort by ▼

## Ormanni 2017 Chianti Classico TUSCANY

Violet, red berry and dark spice aromas lift out of the glass. ...

[SEE FULL REVIEW ▶](#)

**91**  
Points  
\$19

## Tenuta Santa Maria 2013 Riserva (Amarone della Valpolicella Classico) VENETO

This opens with aromas of violet, cooking spices and underbrush. The firm, ...

[SEE FULL REVIEW ▶](#)

**91**  
Points  
\$88

## Contessa 2017 Montepulciano d'Abruzzo CENTRAL ITALY

Dense aromas of brandy-soaked red cherry and plum meet a firm blend ...

[SEE FULL REVIEW ▶](#)

**90**  
Points  
\$19

Violet, red berry and dark spice aromas lift out of the glass. On the savory palate, firm, fine-grained tannins accompany ripe black cherry, blood orange and licorice. Drink through 2022. —*KERIN O'KEEFE*

### RATING

**91**  
POINTS

### PRICE

\$19, [BUY NOW](#)

### VARIETY

[Red Blends](#), Red Blends

### APPELLATION

[Chianti Classico](#), [Tuscany](#), [Italy](#)

### WINERY

Ormanni

[Print a Shelf Talker Label](#)

### ALCOHOL

14%

### BOTTLE SIZE

750 ml

### CATEGORY

Red

### IMPORTER

Virtuoso Selections

### DATE PUBLISHED

3/1/2021

### USER AVG RATING

Not rated yet [[Add Your Review](#)]



# EDA

	name	review	taster	rating	price	designation	variety	appellation	winery	alc	...
0	Fattoria Giuseppe Savini 2017 Rondineto (Monte...)	ripe cherry and wild berry aromas take on a bi...	ALEXANDER PEARTREE	86	12.0	Rondineto	Montepulciano, Italian Red	Montepulciano d'Abruzzo, Central Italy, Italy	Fattoria Giuseppe Savini	13.0	...
1	Feudo Antico 2017 Organic (Montepulciano d'Ab...	spiced raspberry and cherry meld with a bit of...	ALEXANDER PEARTREE	86	10.0	Organic	Montepulciano, Italian Red	Montepulciano d'Abruzzo, Central Italy, Italy	Feudo Antico	13.0	...

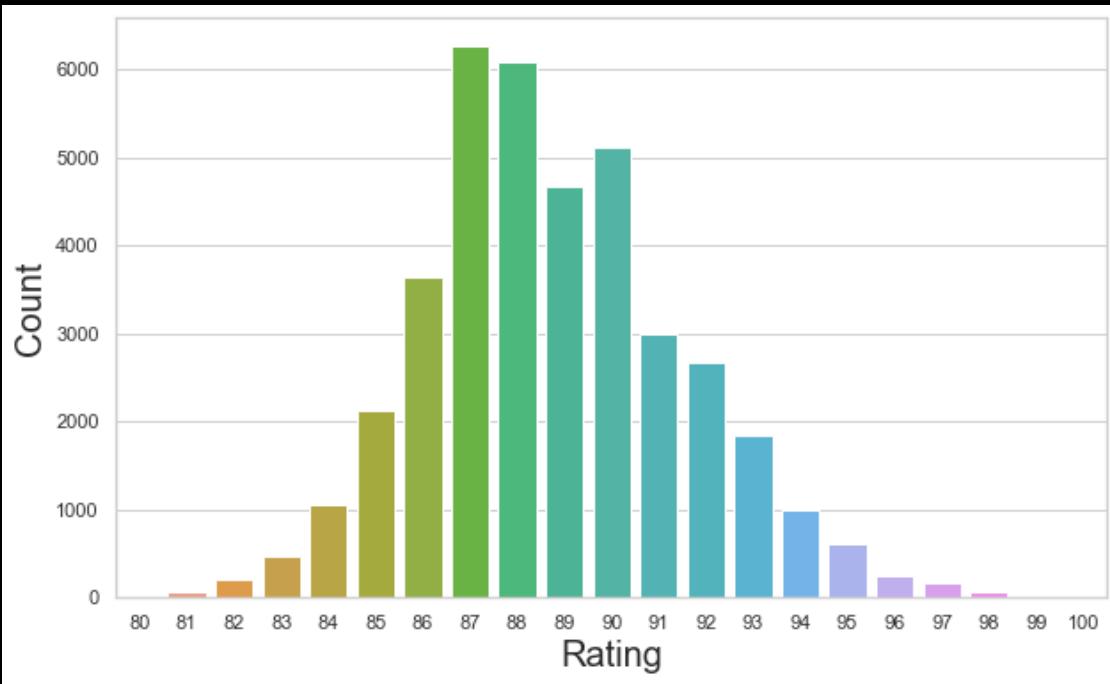
Columns(['name', 'review', 'taster', 'rating', 'price', 'designation', 'variety', 'appellation', 'winery', 'alc', 'bottle\_size', 'category', 'importer', 'review\_date', 'variety1', 'variety2', 'appellation1', 'region', 'nation', 'vintage', 'value\_for\_money', 'final\_judgement', 'final\_score', 'above\_below\_median\_rating', 'review\_test'],

'leather, tilled earth, grated clove and black cherry aromas take center stage along with a whiff of fragrant blue flower. the taut, velvety palate doles out raspberry extract, ripe cranberry, tobacco and licorice alongside firm acidity and close-grained tannins that leave a drying finish. drink

# EDA

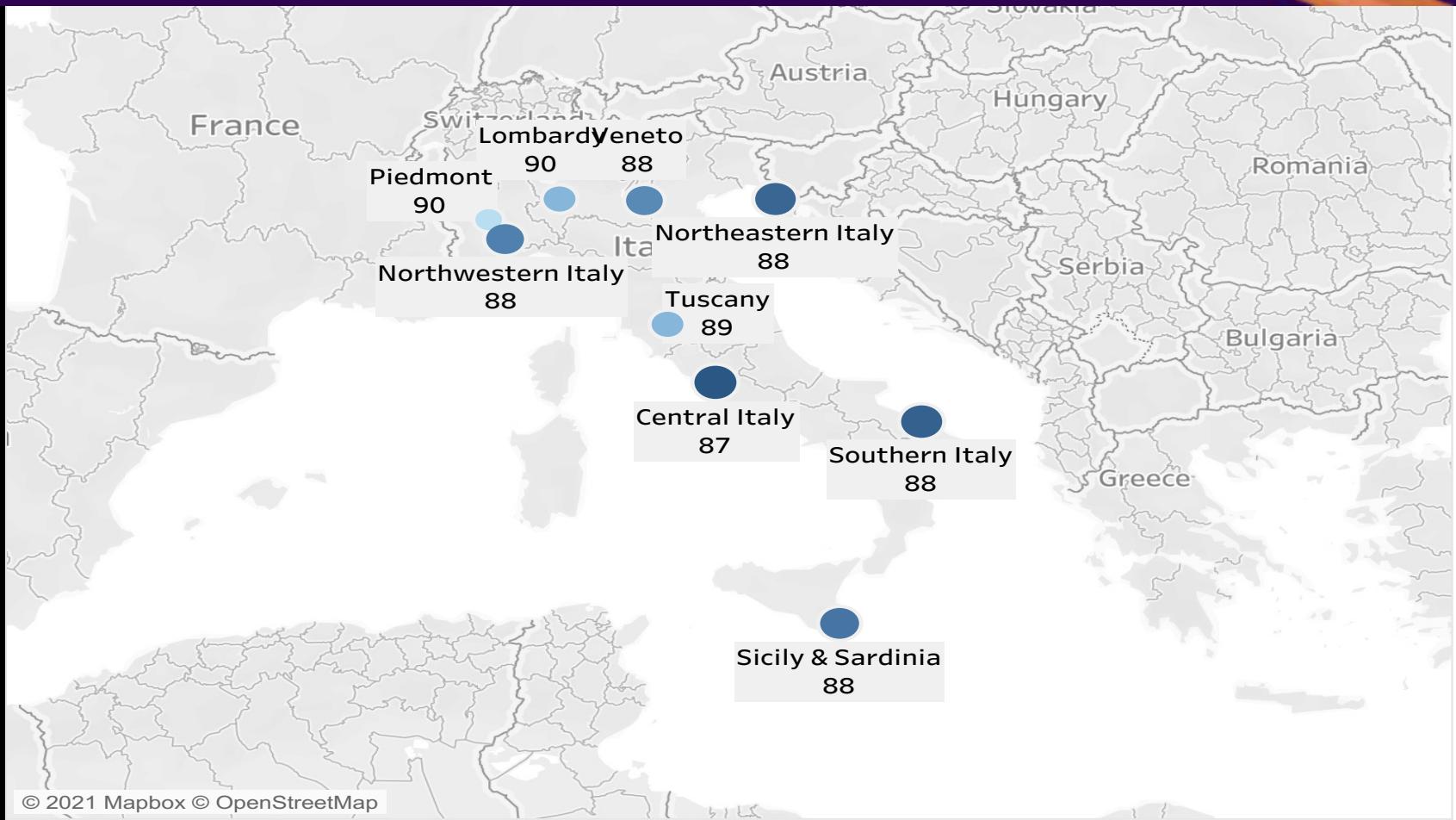
- Wine Enthusiast's 100 - point wine scoring scale:

- 98–100 – Classic.
- 94–97 – Superb.
- 90–93 – Excellent.
- 87–89 – Very good.
- 83–86 – Good.
- 80–82 – Acceptable.



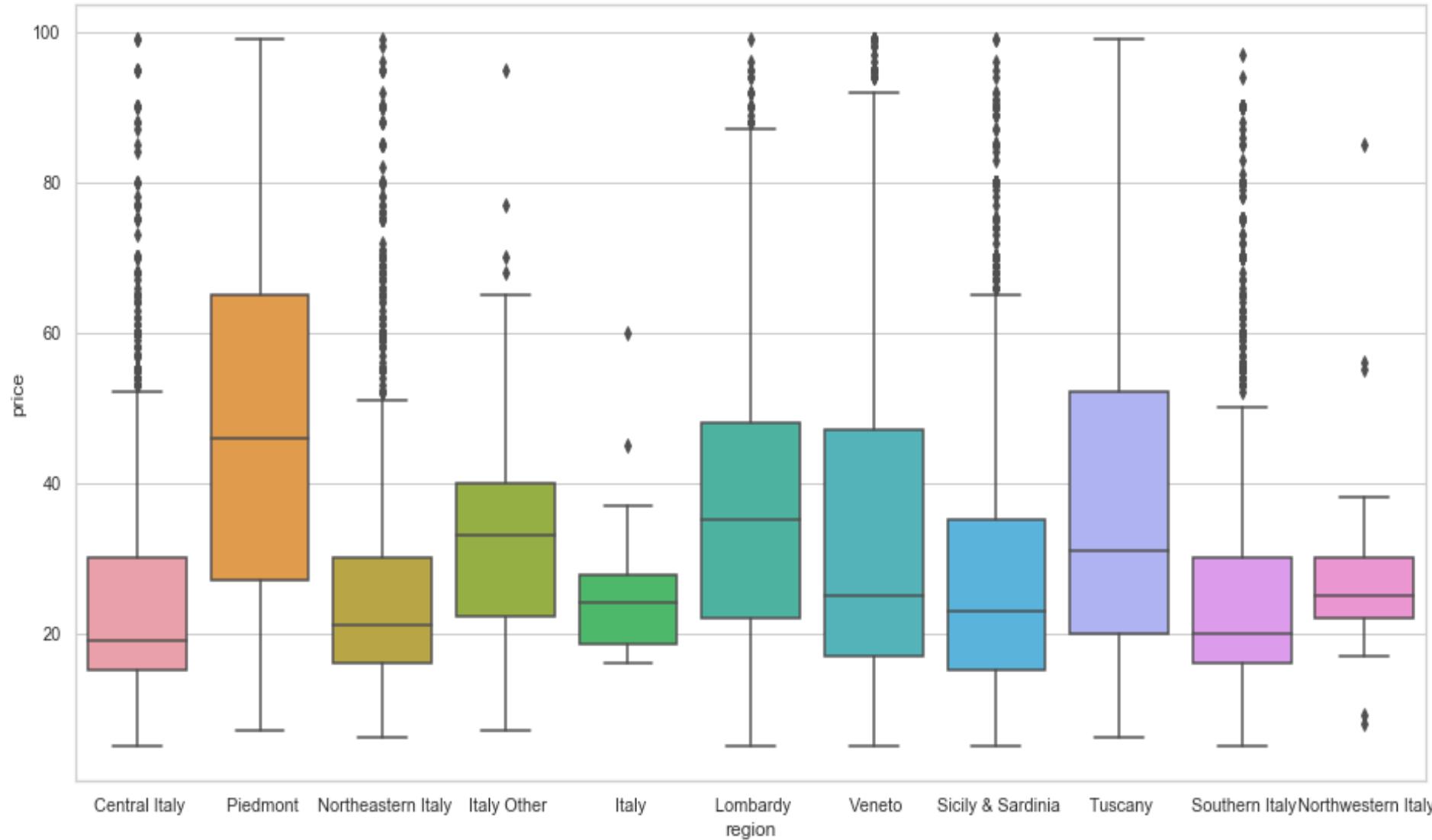


# EDA

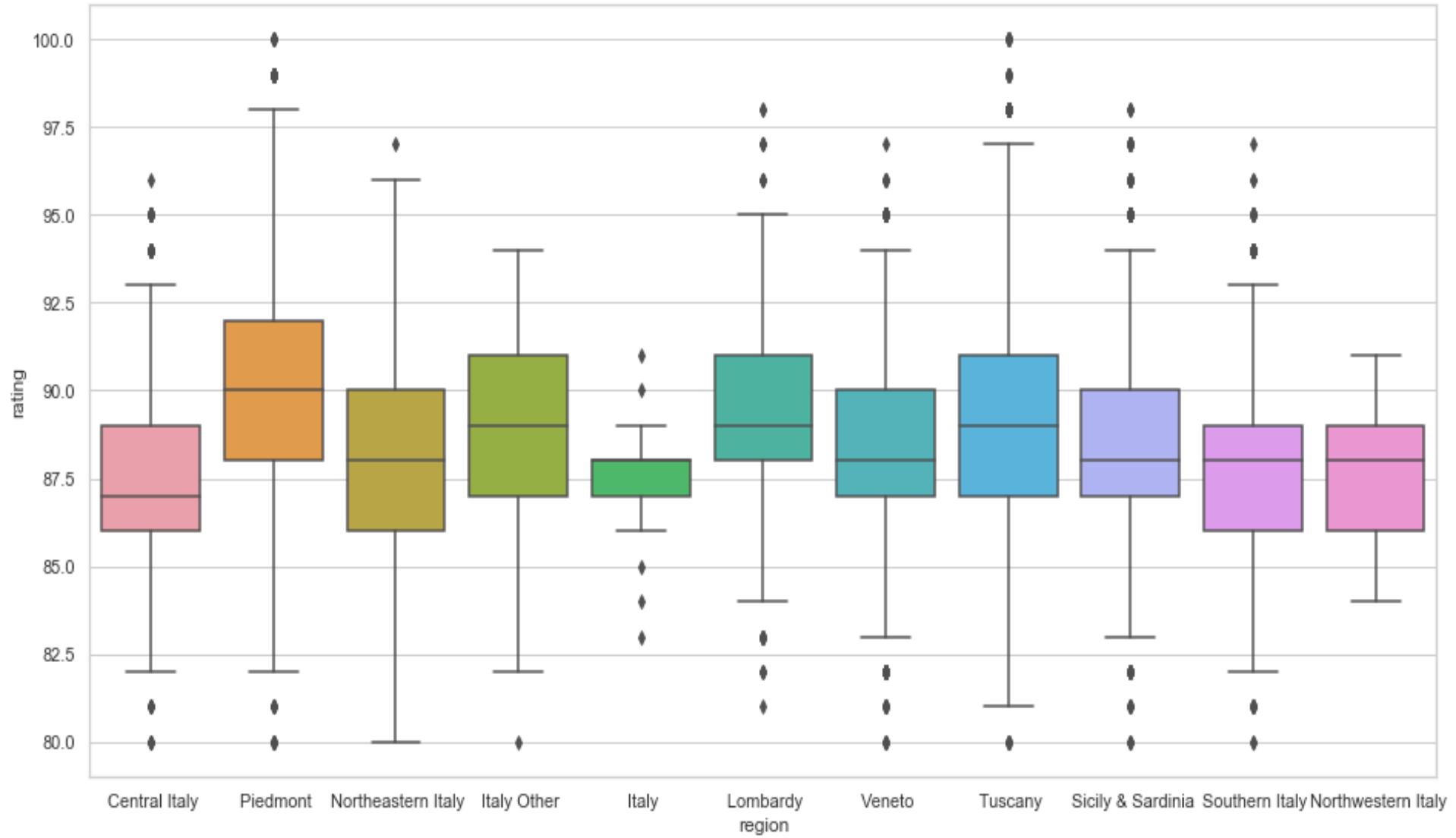


Map based on Longitude (generated) and Latitude (generated). Colour shows median of Value For Money. Size shows median of Value For Money. The marks are labelled by Region and median of Rating. Details are shown for Region. The data is filtered on Value For Money, which excludes Null.

# EDA



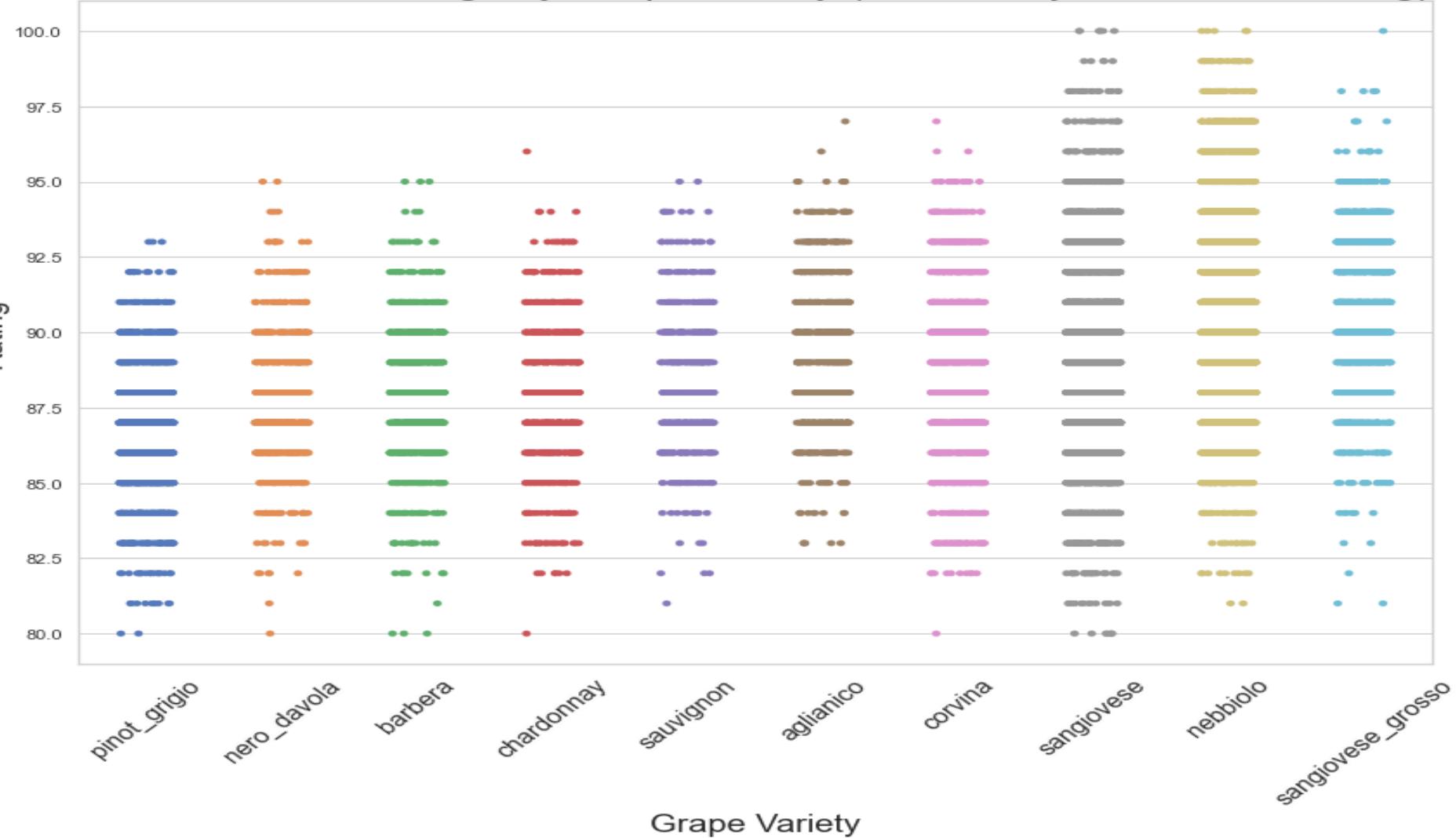
# EDA





# EDA

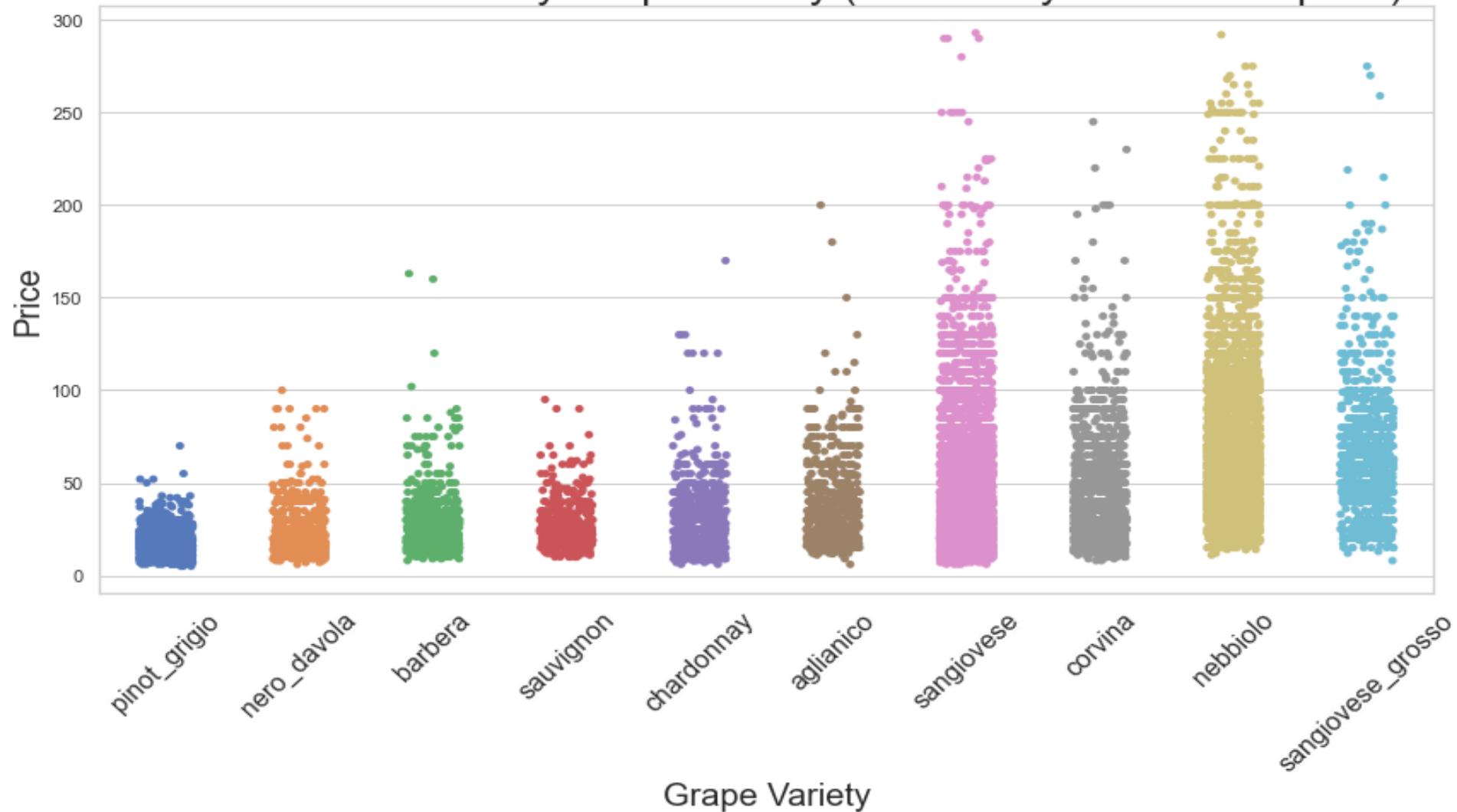
Distribution of Ratings by Grape Variety (ordered by the median rating)





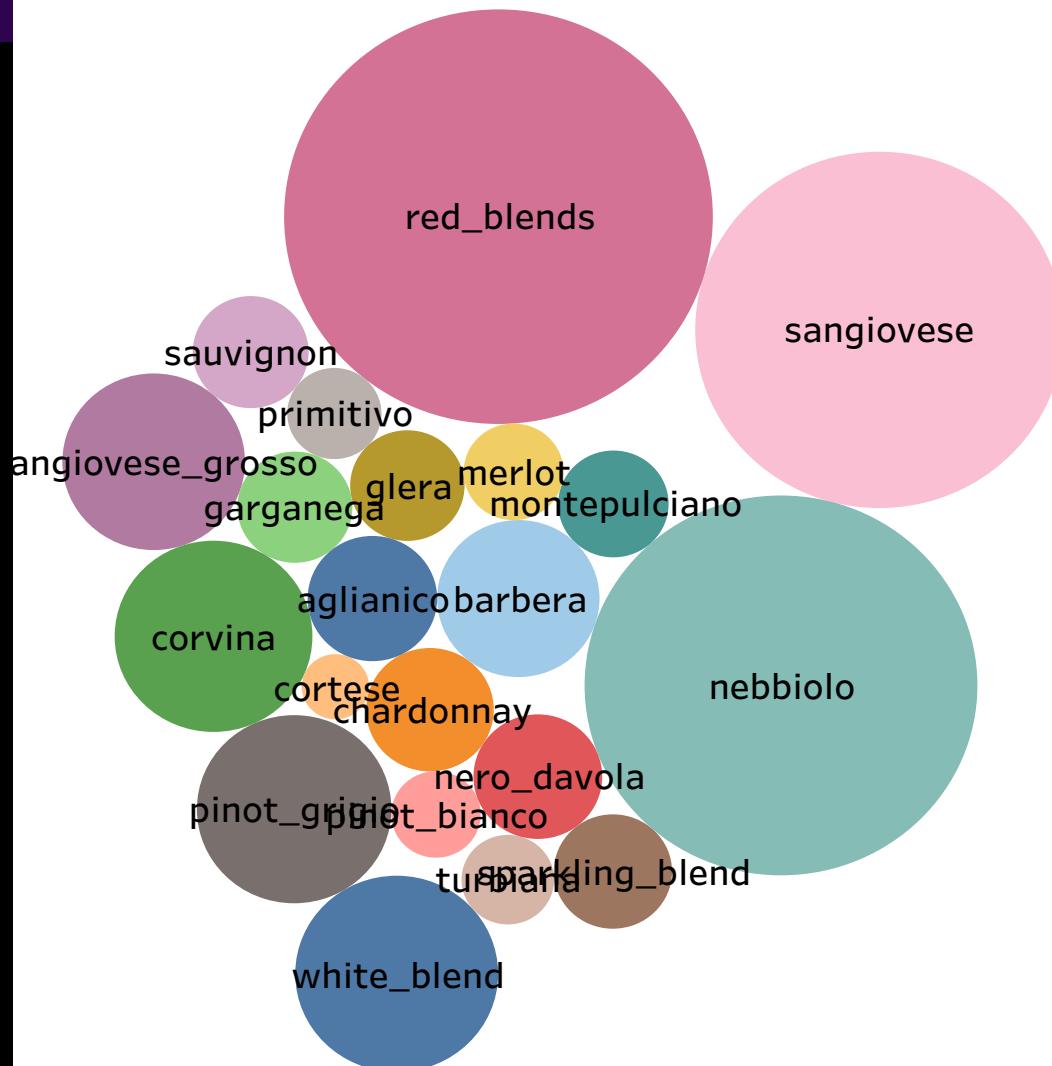
# EDA

Distribution of Price by Grape Variety (ordered by the median price)



# EDA

GrapeVarietySize

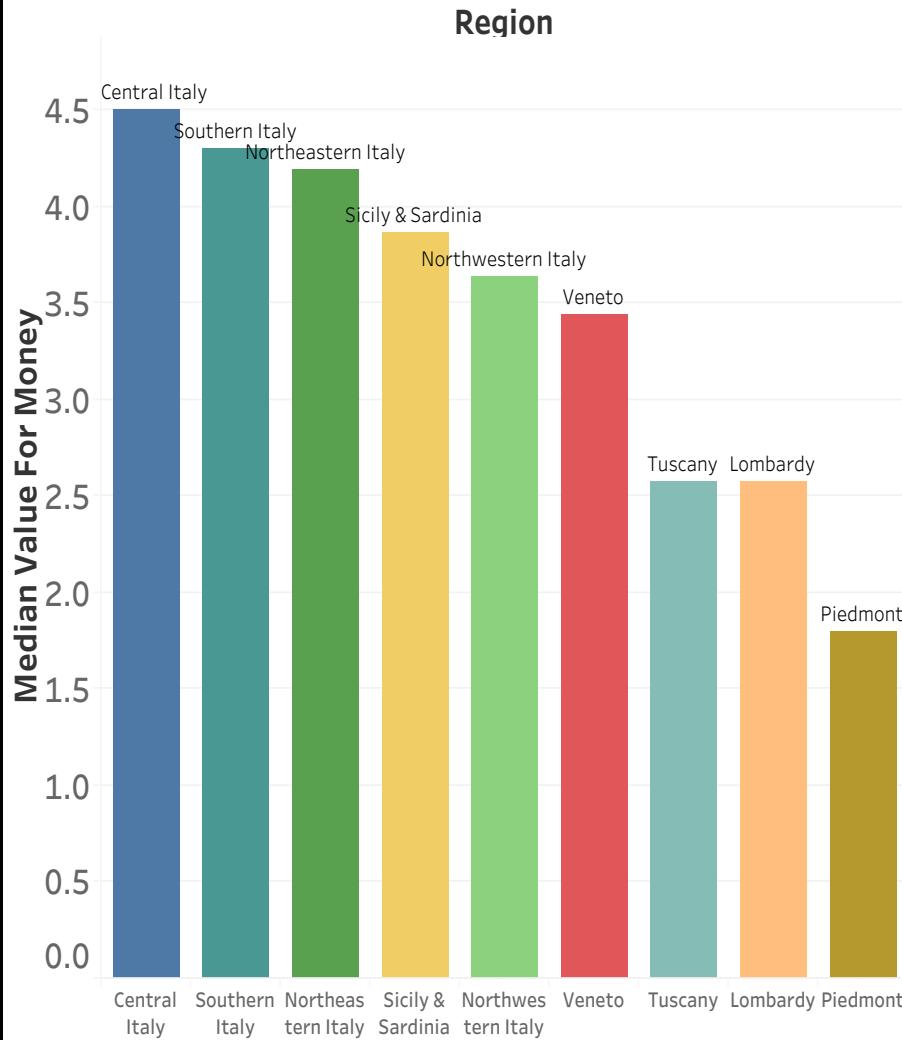


red.blends	7202
nebbiolo	6046
sangiovese	5319
white_blend	1613
corvina	1532
pinot_grigio	1483
sangiovese_grosso	1307
barbera	1033
aglianico	657
nero_davola	650
chardonnay	638
sparkling_blend	548
sauvignon	526
garganega	518
glera	513
montepulciano	478
merlot	391
primitivo	348
turbiana	336
pinot_bianco	313

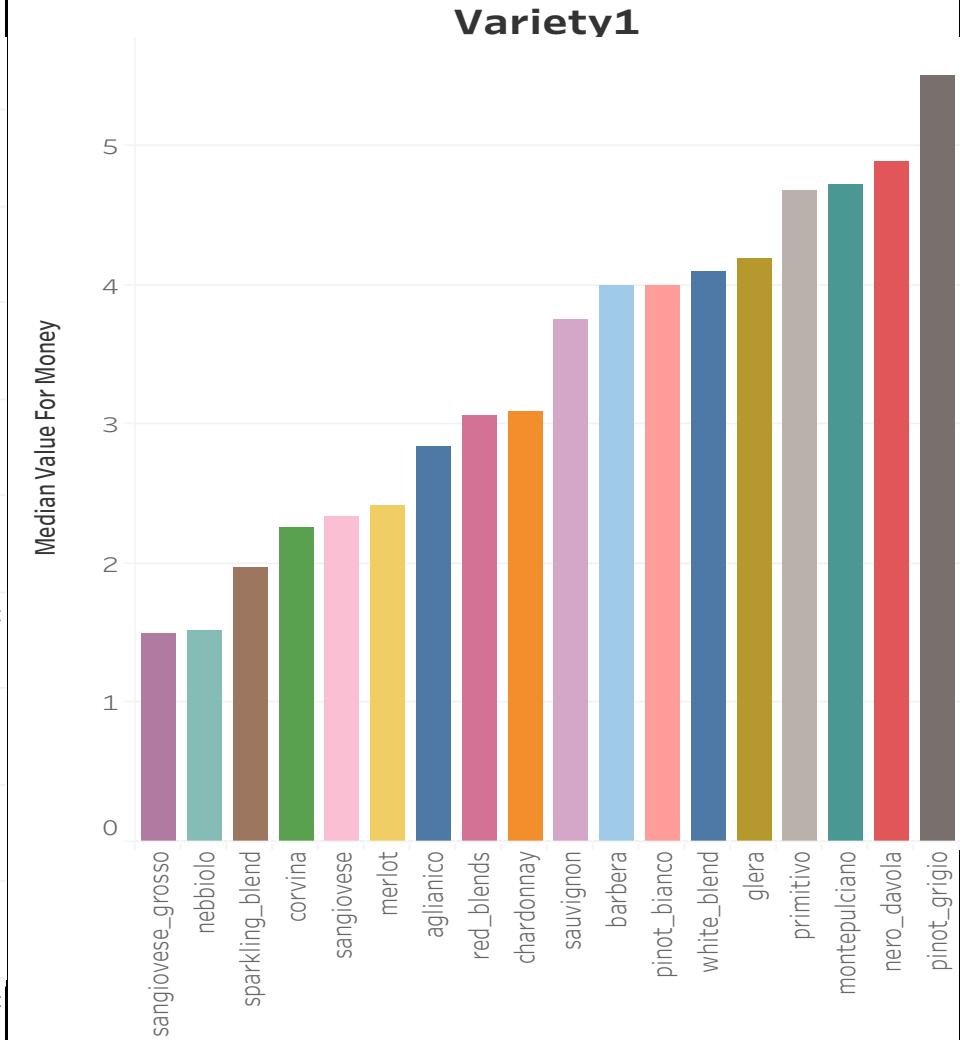
Variety1. Colour shows details about Variety1. Size shows count of Variety1. The marks are labelled by Variety1. The view is filtered on Variety1, which keeps 21 of 231 members.

# EDA

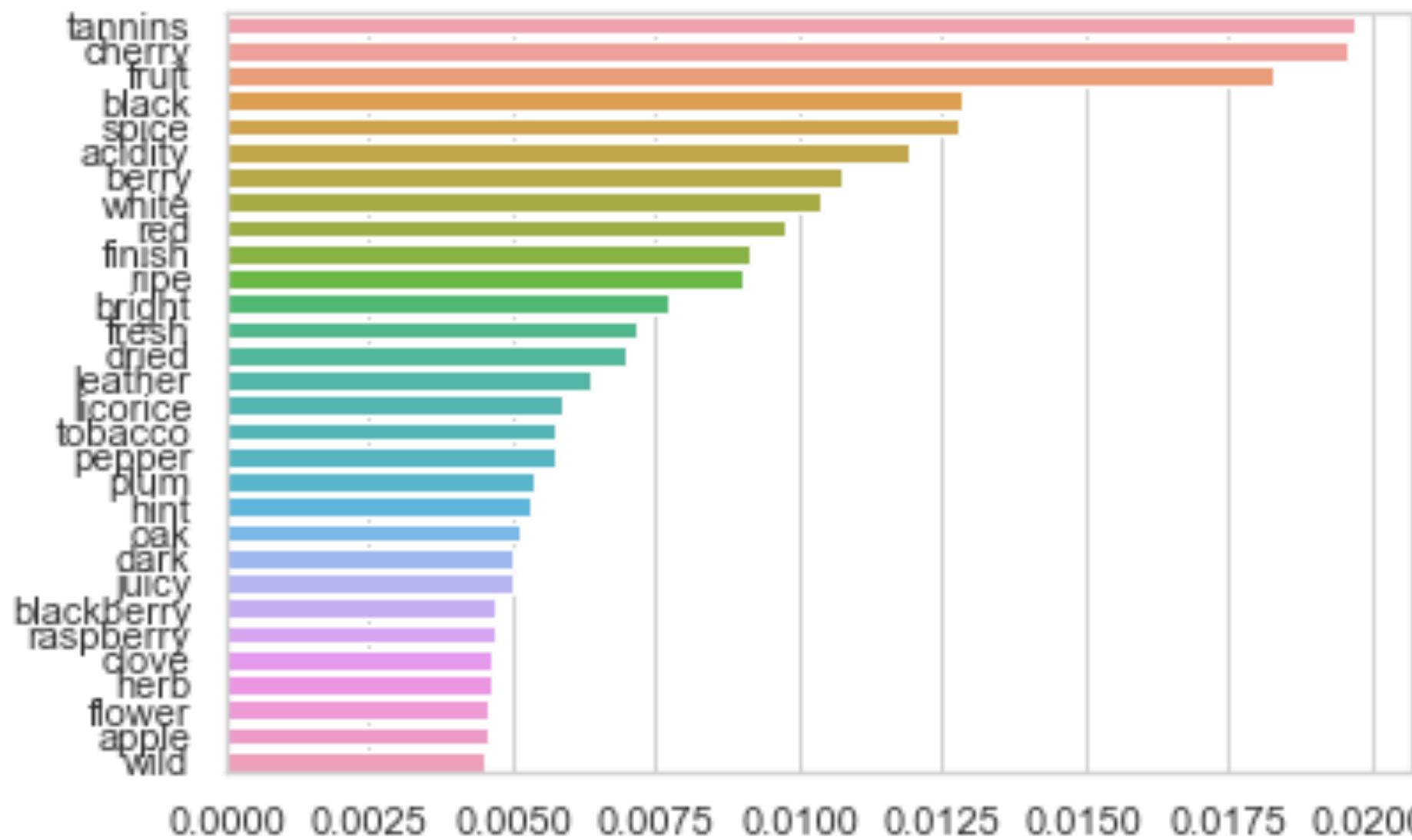
Value\_for\_money and Region



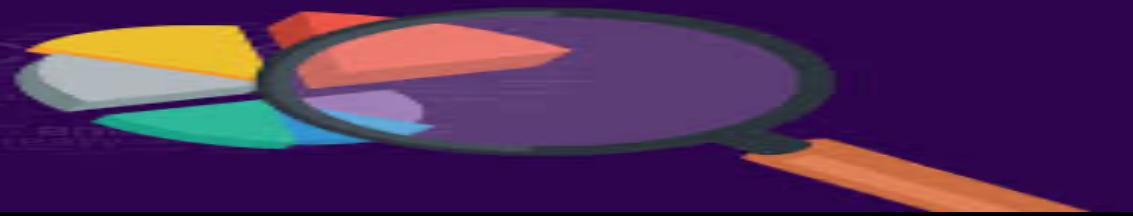
Value\_for\_money and GV



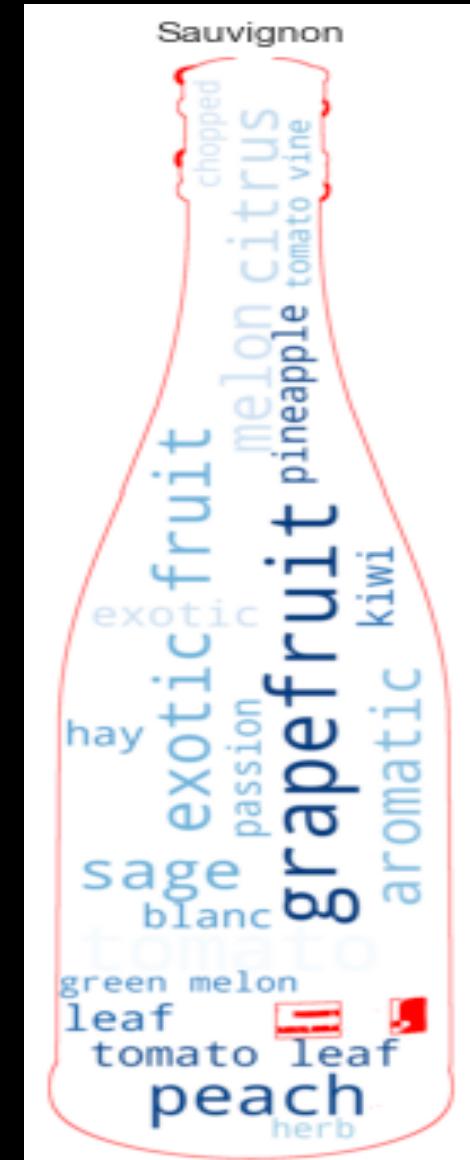
## % Occurrence of Most Frequent Words



# Predict Grape Variety from Reviews – SAUVIGNON BLANC



aromatic  
tomato leaf  
exotic fruit  
exotic sage  
tomato kiwi  
leaf citrus hay  
grapefruit pineapple  
melon peach



# Predict Grape Variety from Reviews



## What's inside the bottle?

Sangiovese

syrah  
chianti classico  
savory herb  
tannins  
tuscany  
wild cherry  
canaiolo  
cherry  
riserva  
plum  
underbrush  
brunello  
cabernet  
tuscan  
mediterranean  
morellino  
thyme  
**chianti**  
central  
classico

Nero d'Avola

**sicilian**  
red  
blueberry  
juicy blackberry  
fruity  
made  
pure expression  
almond  
pistachio  
**sicily**  
tannins soon  
mediterranean  
carob  
**blackberry**  
fleeting  
cocoa  
tannins soft  
imported  
easygoing

Nebbiolo

vineyard  
**barbaresco**  
strawberry  
licorice  
sawdust  
raspberry  
tar  
tannic rose  
truffle  
**barolo**  
already  
cru  
hazelnut  
assertive  
camphor  
**roero**  
tannins  
approachable

Sauvignon

chopped  
melon citrus  
pineapple tomato vine  
kiwi  
exotic fruit  
passion fruit  
aromatic grapefruit  
sage blanc  
blanc  
green melon leaf  
leaf tomato leaf peach  
herb

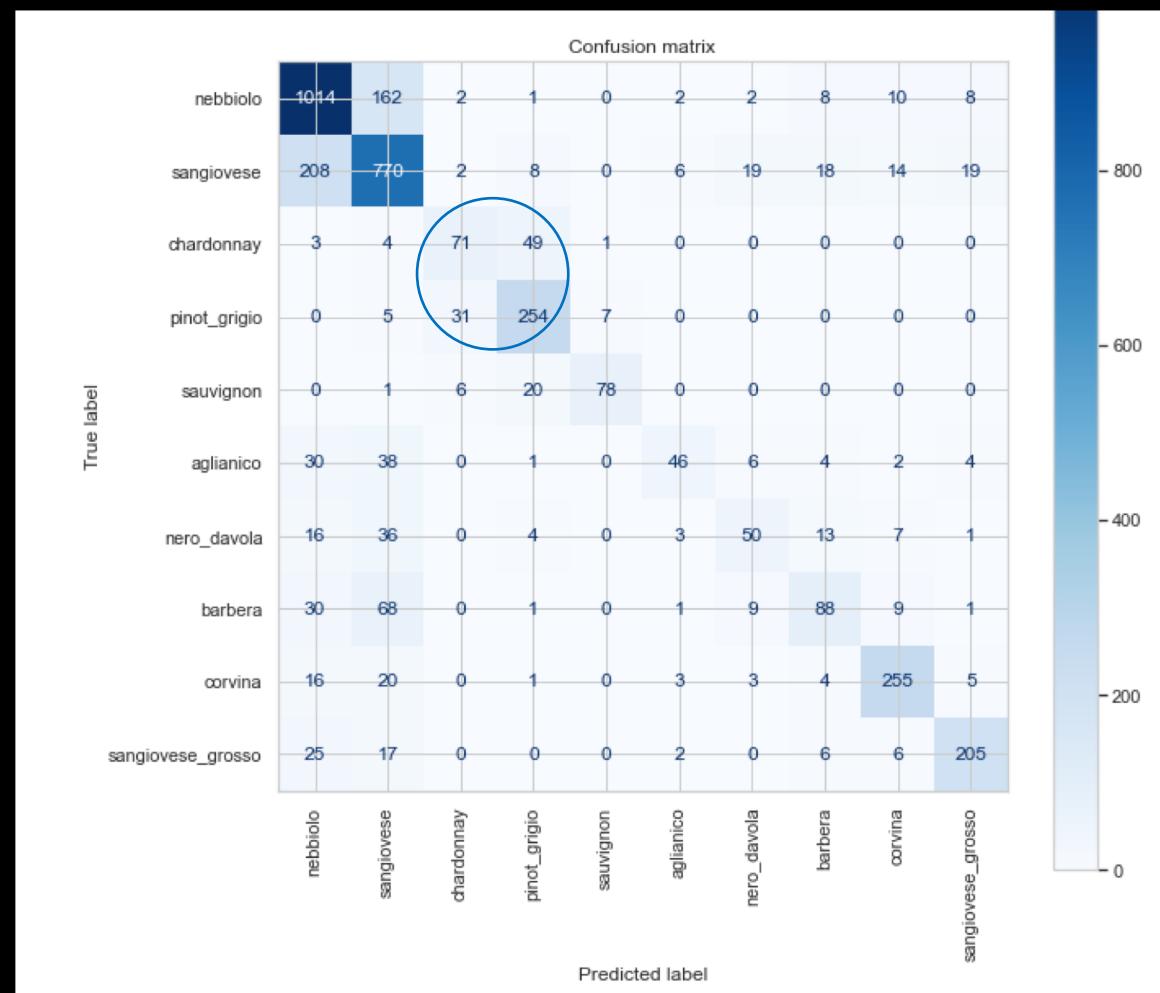
# Predict Grape Variety from Reviews

Cvec-Tvec Lr

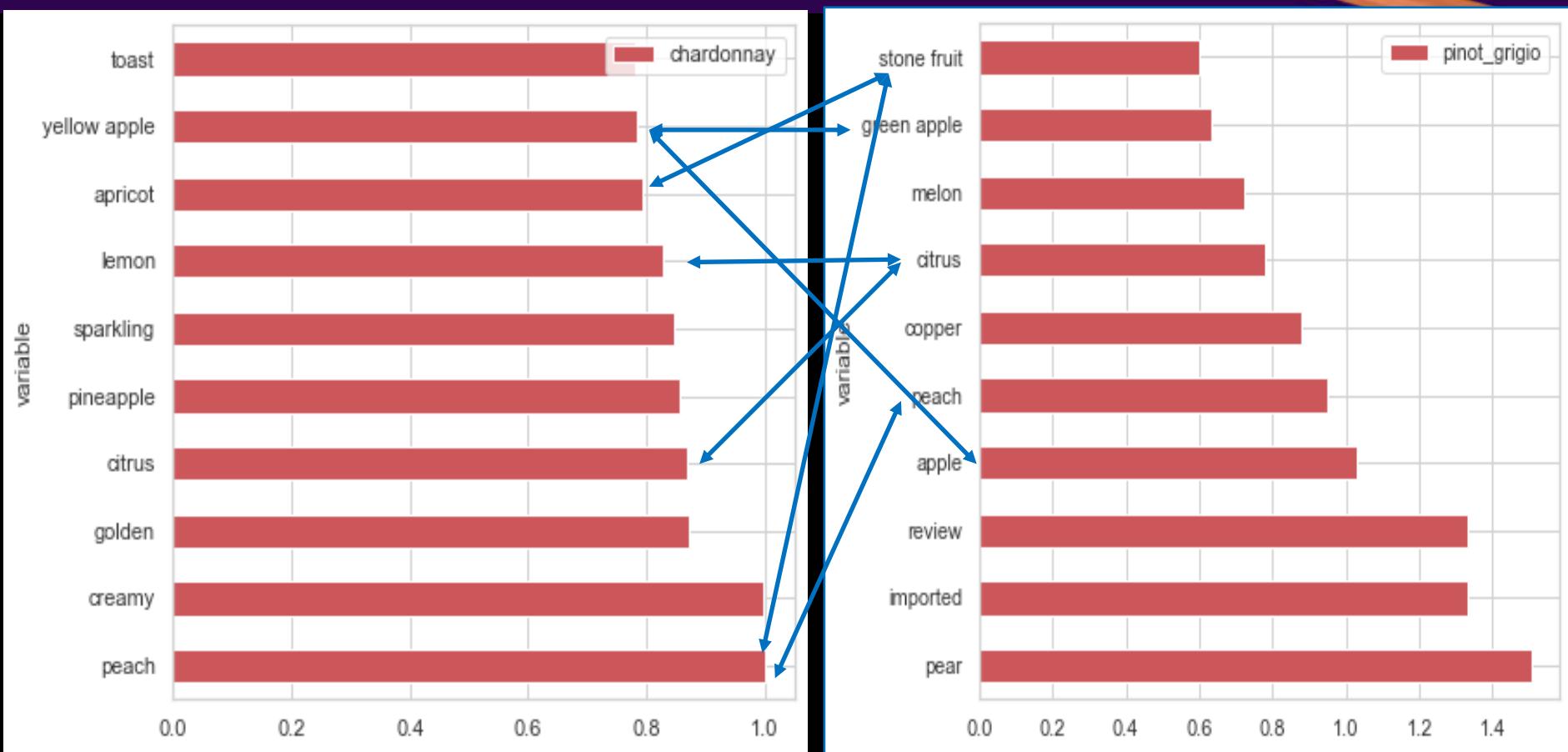
CV 0.739

SVM 0.7

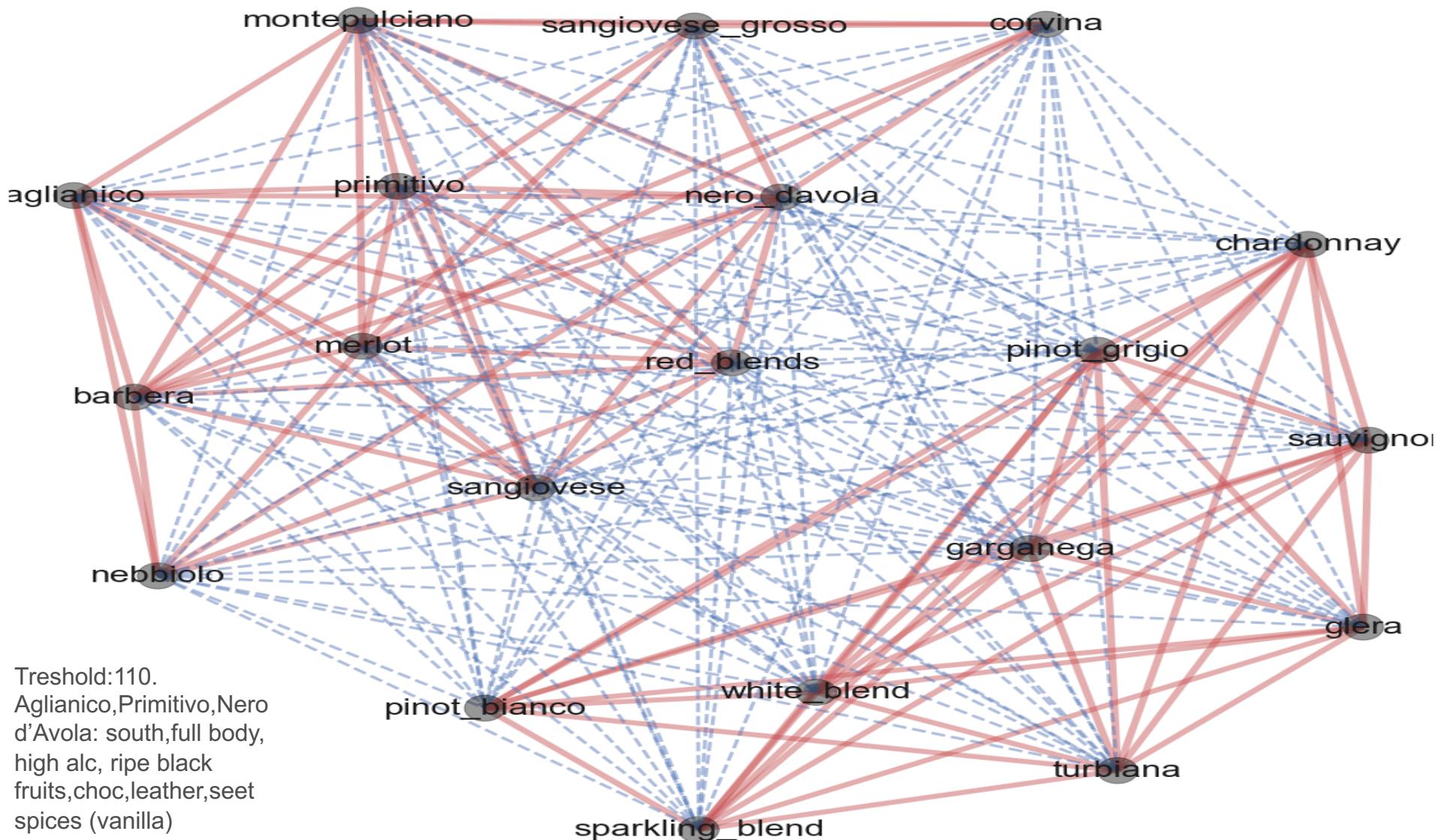
	precision	recall	f1-score
aglianico	0.73	0.35	0.47
barbera	0.62	0.43	0.51
chardonnay	0.63	0.55	0.59
corvina	0.84	0.83	0.84
nebbiolo	0.76	0.84	0.79
nero_davola	0.56	0.38	0.46
pinot_grigio	0.75	0.86	0.80
sangiovese	0.69	0.72	0.70
sangiovese_grosso	0.84	0.79	0.81
sauvignon	0.91	0.74	0.82
accuracy			0.74
macro avg	0.73	0.65	0.68
weighted avg	0.73	0.74	0.73



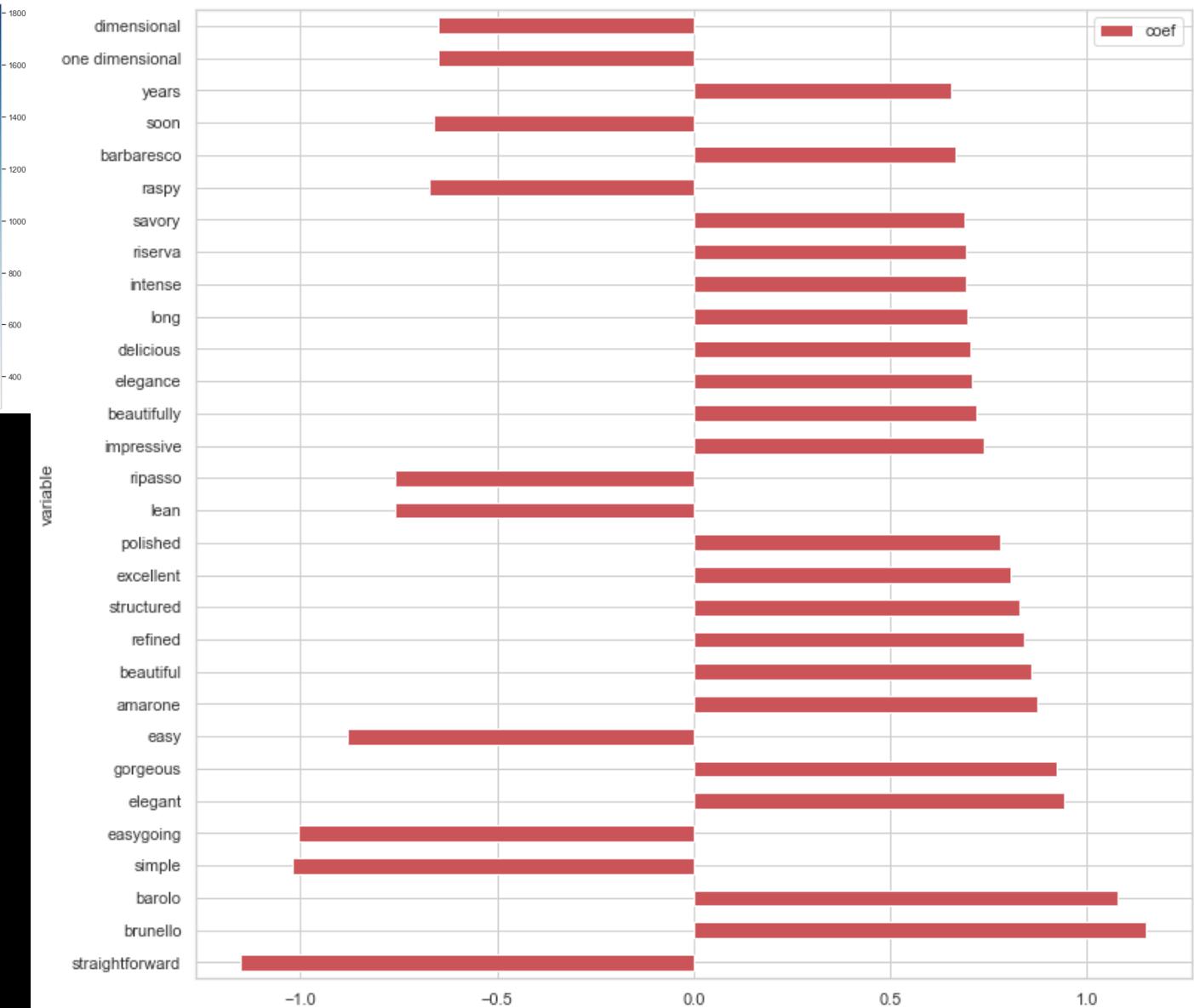
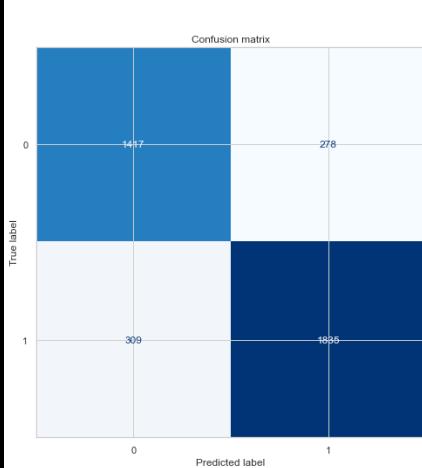
# Predict Grape Variety from Reviews



# How are the grape varieties connected in terms of common flavor characteristics?



# Predict above or below median rating from the text reviews

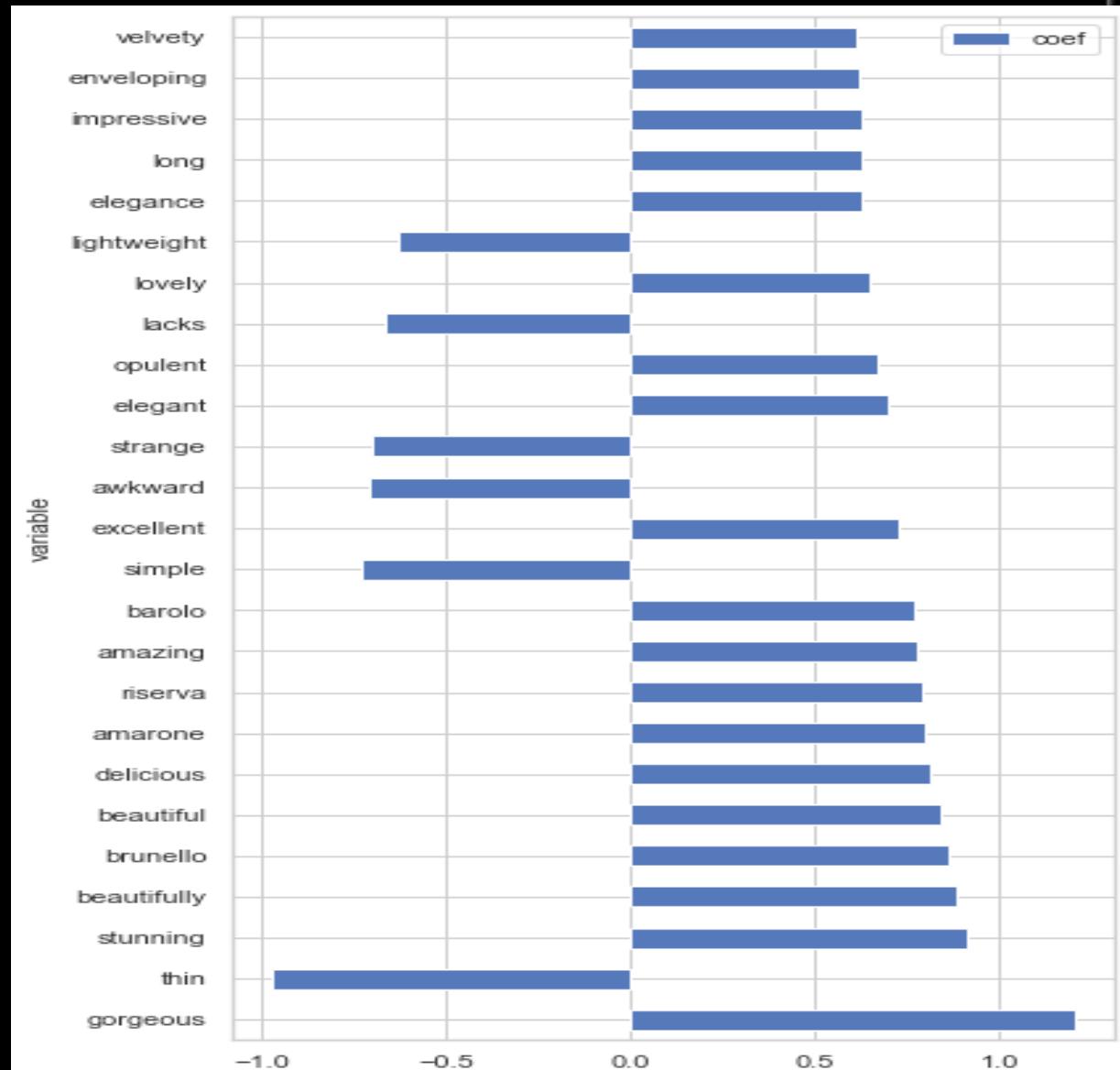


LR CV= 0.84 (train  
.99)

# Predict ratings from the text reviews

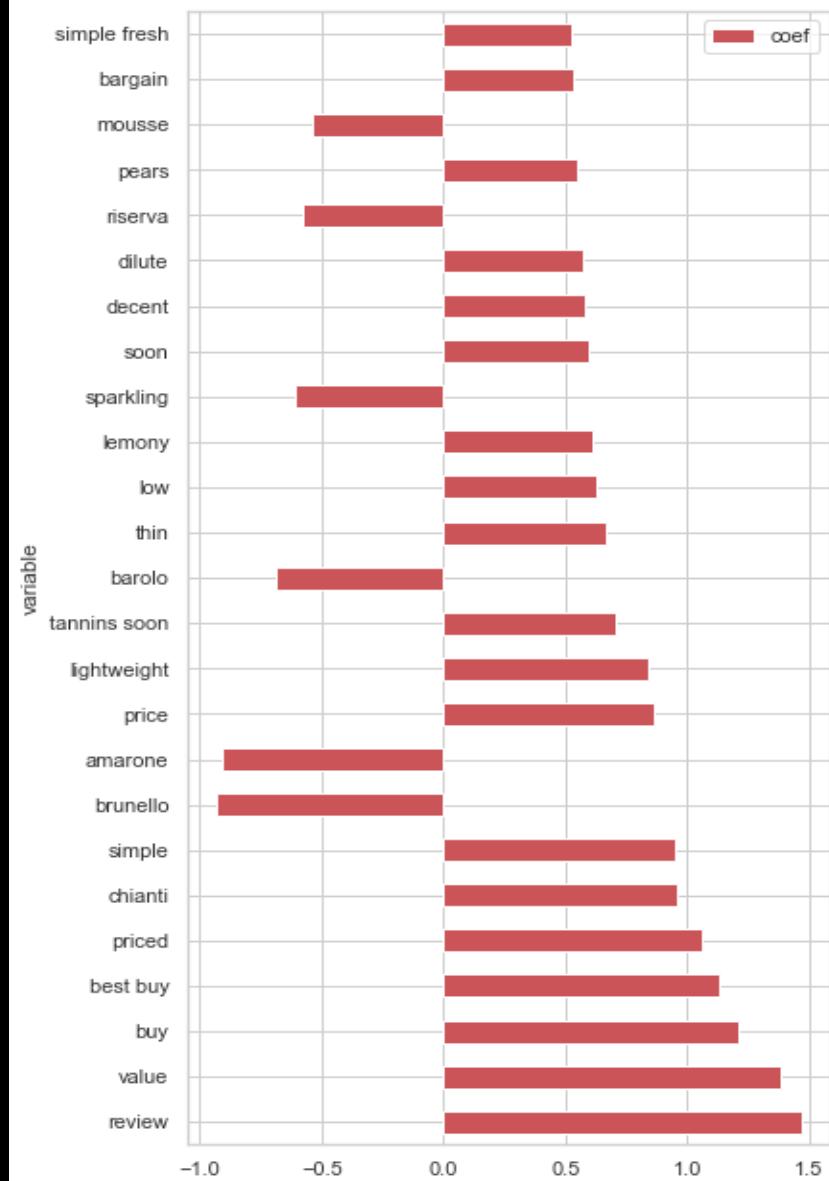
Regression → Ridge Cv = .75  
(92 train score)

Random Forest-Dec Tree CV  
.25



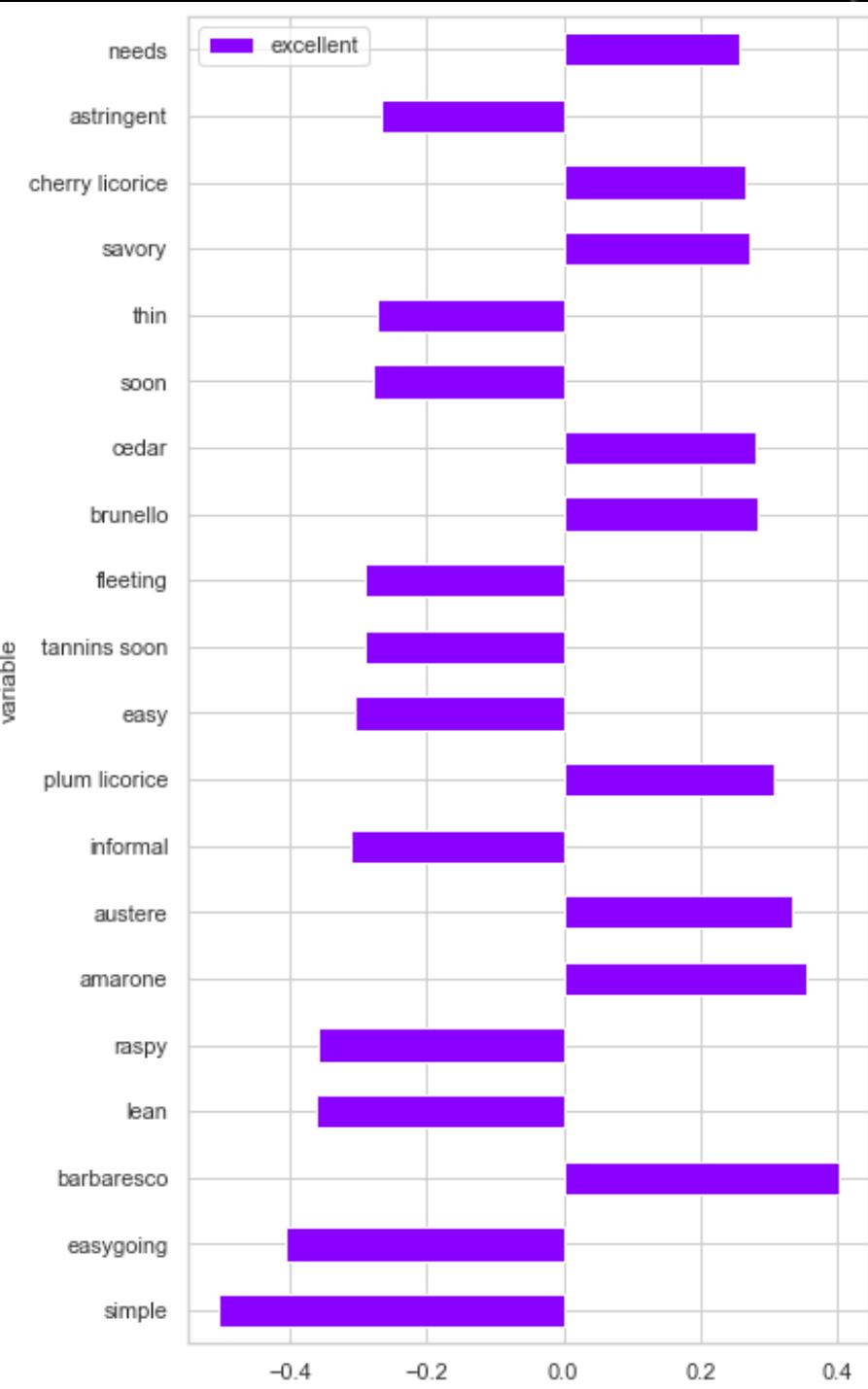
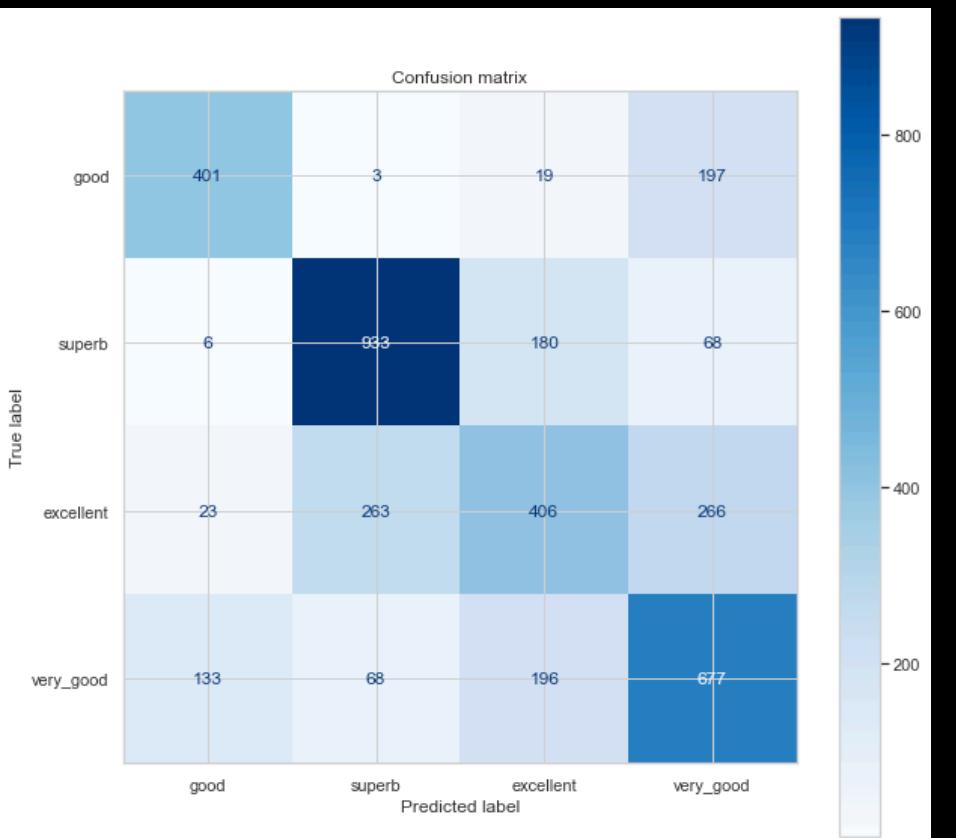
# Predict Value for Money from the text reviews

Regression → Ridge Cv = .54  
(86 train score)  
Lasso/Net much lower



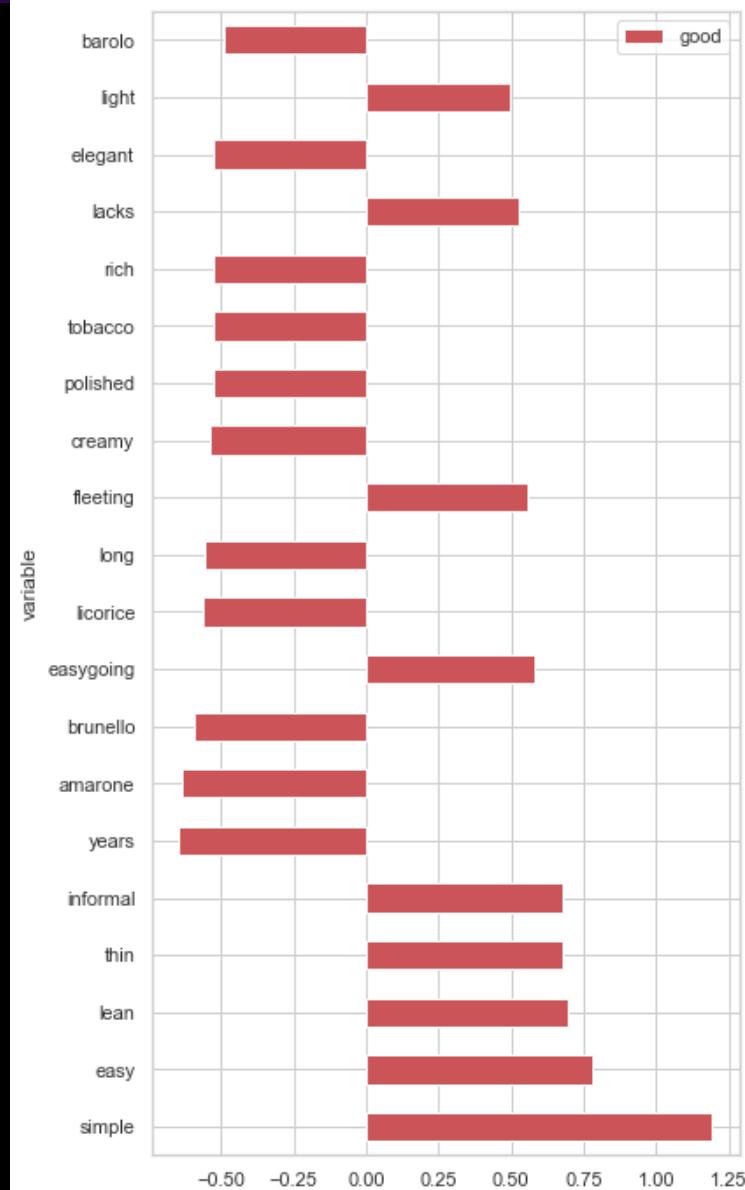
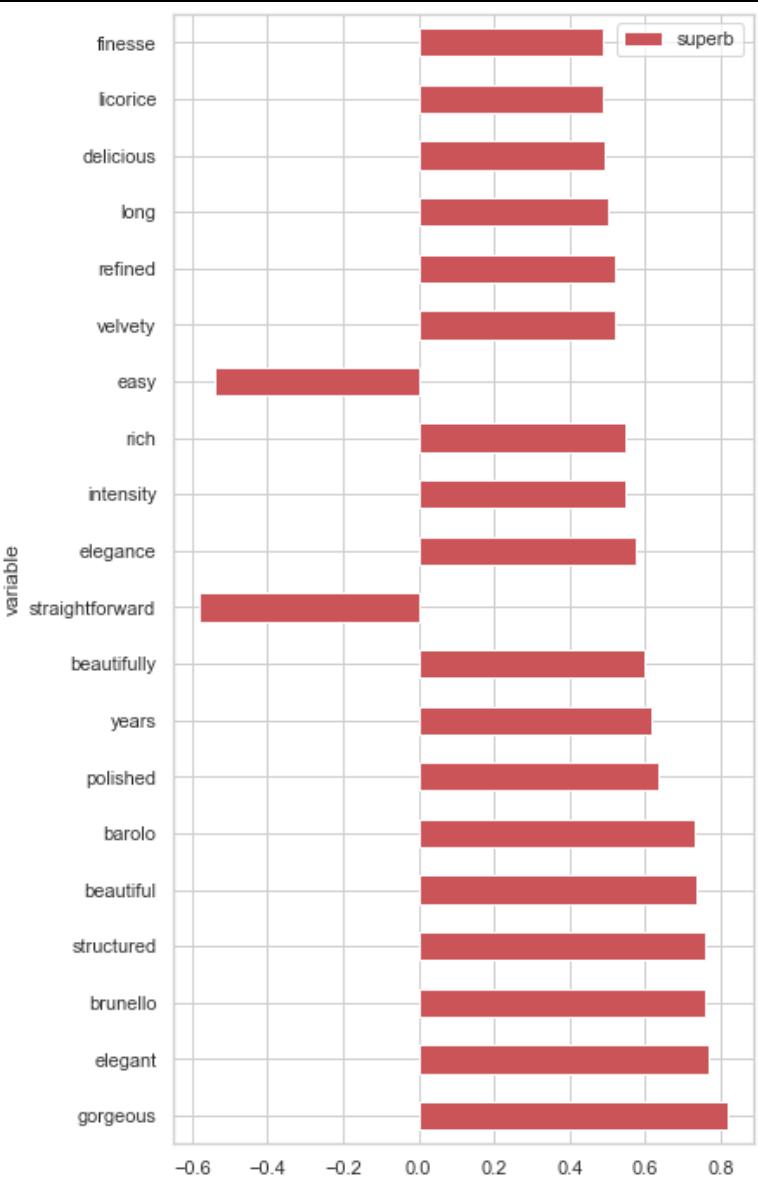
# Predict Good,VeryGood,Excellent, Superb from the text reviews

Classif. → LR CV = .63 (89  
train score)



# Predict Good,VeryGood,Excellent, Superb from the text reviews

Classif. → LR CV = .63 (89  
train score)



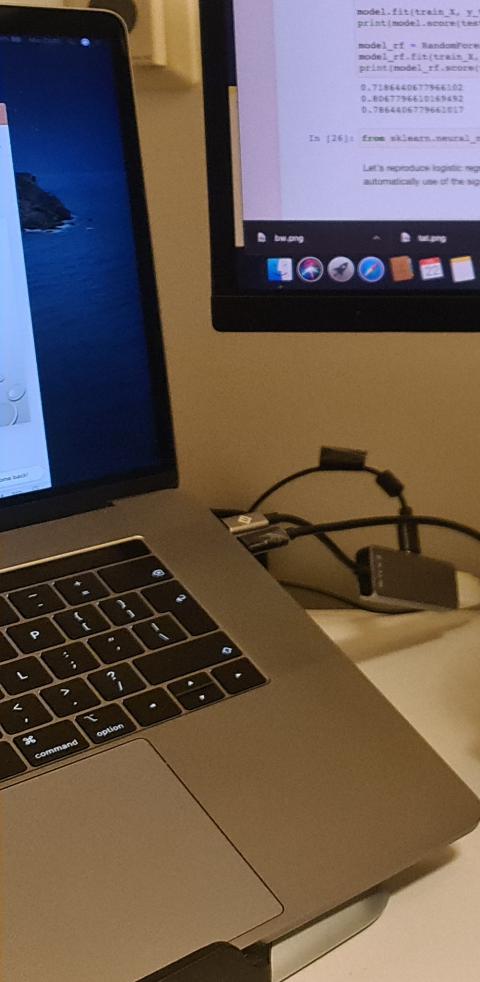
# Predict above/below median Rating from Appellation

Classif. →  
LR CV = .70  
(70 train  
score)



A photograph of a paved road curving away from the viewer towards a bright, overexposed horizon. The sky above the horizon is a vibrant blue. Overlaid across the center of the image is the question "WHAT'S NEXT?" in large, white, sans-serif capital letters.

WHAT'S NEXT?



A large bottle of Le Volte dell'Ornellaia 2015 wine stands prominently in the foreground. The label features a crown at the top, followed by the text "LE VOLTE DELL'ORNELLAIA" and "2015". Below this is a small landscape photograph of a vineyard. At the bottom of the label, it says "TOSCANA" and "INDICAZIONE GEOGRAFICA TIPICA".

Mon\_22\_feb\_evening Last Checkpoint: 17 minutes ago (autosaved)

In [23]: train\_X.shape, test\_X.shape  
Out[23]: ((596, 9), (295, 9))

In [24]: from sklearn.linear\_model import LogisticRegression, Perceptron  
from sklearn.ensemble import RandomForestClassifier

In [25]: model = Perceptron(tol=10\*\*(-6))  
model.fit(train\_X, y\_train)  
print(model.score(test\_X, y\_test))

model = LogisticRegression(fit\_intercept=True,  
C=10\*\*10,  
solver='lbfgs',  
random\_state=1,  
max\_iter=1000)  
model.fit(train\_X, y\_train)  
print(model.score(test\_X, y\_test))

model\_rf = RandomForestClassifier(n\_estimators=500)  
model\_rf.fit(train\_X, y\_train)  
print(model\_rf.score(test\_X, y\_test))

0.7186440677964102  
0.806779641019482  
0.788440677966202

In [26]: from sklearn.neural\_network import MLPClassifier

Let's reproduce logistic regression. The hidden layer contains the identity function whereas the automatically use of the sigmoid.

In [27]: X\_train, X\_test, y\_train, y\_test = train\_test\_split(X\_rw, y\_rw, test\_size=0.2, random\_state=1)

# fit the model  
model.fit(X\_train, y\_train)  
# evaluate on the training set  
print('Best alpha:', model.alpha\_)  
# evaluate on the training set  
print('Training score:', model.score(X\_train, y\_train))  
# evaluate on the test set  
print('Test Score:', model.score(X\_test, y\_test))

Best alpha: 0.3593813663042626  
Training score: 0.4714463652948817  
Test Score: 0.3462450631109546

In [28]: model.alpha\_

Out[28]: 0.3593813663042626  
# fit the model instance  
(alpha=alpha)  
# fitted scores  
# al\_scores(model, X\_train, y\_train, cv=5)  
# fitted training scores  
# al\_scores(model, X\_train, y\_train, cv=5)  
# fitted training scores  
# al\_scores(model, X\_train, y\_train, cv=5)  
# the data on the whole training set  
# al\_scores(model, X\_train, y\_train, cv=5)  
# fit the test set  
# al\_scores(model, X\_train, y\_train, cv=5)  
# fitted test scores  
# al\_scores(model, X\_train, y\_train, cv=5)

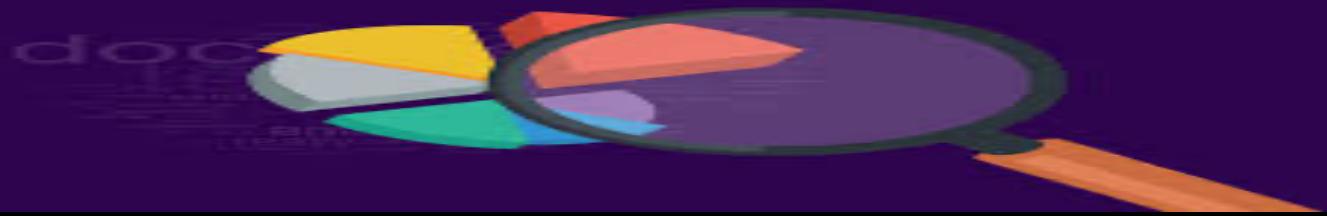
Training scores: [0.34214506 0.34397601 0.33855584 0.3529229 0.35313324]  
Testing scores: [0.34614260940003616 0.34652948827 0.34652948827 0.34652948827]

1109546

low average RATING FROM THE WINERY?



## LIMITS & IMPROVEMENTS



**DATASET** : biased towards red wines, need more data and expand beyond Italy

**MODELS**: more and more

**CLASSIFIER**: try to create one with more data

**CLUSTERING**: more and other techniques



