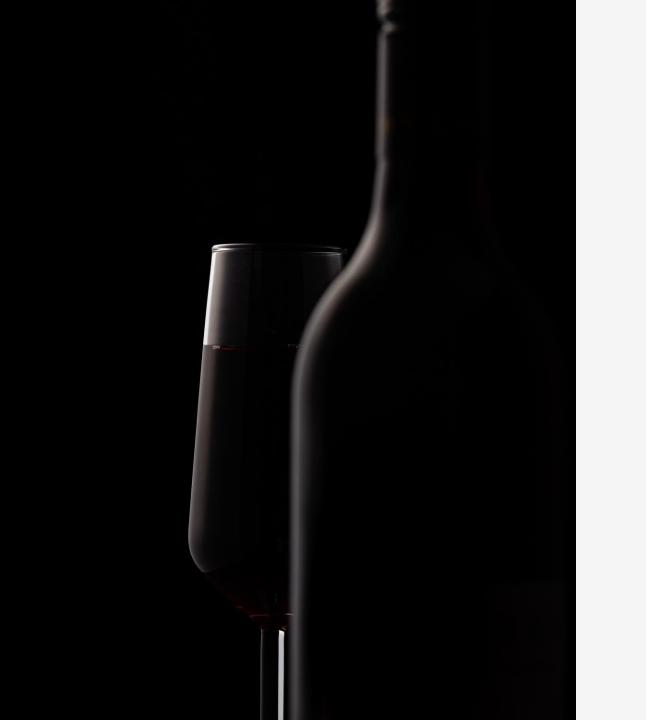


ALL YOU NEED IS...GOOD WINE

3251 - 017.'

Group 6: Ali, Neda, Nisarg, Olga, Vaishali, Vinay



Agenda

- 1. Introduction
- 2. Hypothesis Testing
- 3. Logistic Regression
- 4. Linear Regression and SVM
- 5. Conclusions

HOWMUCH DOWE DRINK?

What does an average month of drinking look like?

10 Days of 13 Unique consumption occasions per month

25 Total drinks consumed

3.6 liters more than the world average

Across multiple categories (beer, wine, spirits, cider, RTD)

MARKET VALUE OF THE WINE INDUSTRY IN CANADA 8.59bn CAD

FORECASTED MARKET VALUE OF THE WINE INDUSTRY IN CANADA

10.85bn CAD



HOWWEDRINK?











· Only on weekends

DRINKS TO

Weekends and

occasional evenings

To relax and have fun

DRINKS WITH

Evenings at home

With her husband

. To relax with a meal

1 or 2 glasses with dinner

MEALS

. In and out of the house

UNWIND

. To connect and have fun

Usually out of the house

DRINKER





- · A variety of craft beers
- · Cider once in a while
- · Spirits later at night
- · Favourite brand is Creemore



- Wine around the home
- · Mixed drinks with friends
- Favourite wine is Jackson Triggs

JUST WINE

- Only wine
- Many varietals and brands under \$20
- · No favourite brand

HIGHBALLS AND WINE

- · Wine with dinner
- · Canadian whisky at cocktail
- · Favourite whisky is Crown Royal





- · Loads up every few weeks
- · For home and to take to a friend's house
- . Decides what to buy while in the store - cold if possible

GO TO **BRAND**

- · Every other week
- To take home for routine drinking
- Buys her go-to brand

VALUE **SEEKER**

- Weekly
- To take home for routine drinking or bring to a friend's
- Looks for in-store specials

BUYS ON DISPLAY

- Weekly
- · Wine to go with dinner
- . Looks for good wine on display



Gen X

Jennifer

Boomer

Brenda

Retired

Richard



URBAN MILLENNIAL

- Single with roommates
- Urban professional
- Works hard and plays hard

SUBURBAN SPOUSE

- Married with two kids
- Suburban professional
- Manages work and family
- Needs to relax



- · Married, no kids at home
- Urban empty nester
- Many interests and hobbies
- Reads print and online





- Retired in cottage country
- Tries to stay busy
- · Watches a lot of TV







ROUTINE DRINKER

- Evening routine
- To relax and snack
- Usually at home
- · Alone or with his wife

HOW WE DECIDE WHAT TO DRINK?



BUT WHAT ABOUT QUALITY?

Let's consider hypotheses that can help us understand the quality of wine

DATA DESCRIPTION



Records

• White: 4,898

• Red: 1,599

Data types

Real, integer & string

12 attributes, incl. quality ranking
No missing values

Acids

Fixed, Volatile - high levels can lead to a vinegar taste Citric acid - in small quantities adds freshness and flavor

• Sulfur

Free SO2, Sulphates & Total SO2 - antioxidant and antimicrobial

Alcohol

Ethanol % by volume

<u>Residual sugar</u>: amount of sugar remaining after fermentation Chlorides: sodium chloride

<u>Density</u>: will differ from water (1) depending on the alcohol and sugar content <u>pH</u>: 0 (very acidic) to 14 (very basic); most wines are between 3-4 on the pH scale <u>Quality</u>: score between 0 and 10 (rated by sommeliers)

Source of data: https://archive.ics.uci.edu/ml/datasets/Wine+Quality



HYPOTHESIS TESTING

HYPOTHESIS TESTING

HYPOTHESIS 1



HYPOTHESIS 2



HYPOTHESIS 3



NULL HYPOTHESIS

Mean quality scores for red and white wines are the same

Red and white wines have the same average alcohol levels

Quality is the same for wines with higher and lower than average abv%

ALTERNATIVE HYPOTHESIS

Mean quality scores for red and white wines are different

There is difference in alcohol levels between red and white wines

Wines with higher ABV% are of significantly better quality than wines with lower ABV%

OUTCOME

MEAN QUALITY WINE IS HIGHER FOR WHITE WINE THAN IT IS FOR RED WINE

MEAN ALCOHOL LEVEL IS HIGHER FOR WHITE WINE THAN IT IS FOR RED WINE

WINES WITH HIGHER ALCOHOL CONTENT ARE PREFERRED

EXPLANATION

The region is known for its white wines (86%) and is the main source of a unique style of white wine. Portugal has the advantage of numerous indigenous grape varieties that are able to keep their acidity in hot climates and including these in a blend provides balancing freshness for rich white wines.

Given the climate in the North-West corner of Portugal, the grape varieties that grow in that region make wines with high ABV%.

Alcohol content of wine has spiked considerably in recent years. Demand is high for intense flavors and higher alcohol content.

TOTALS







LOGISTIC REGRESSION

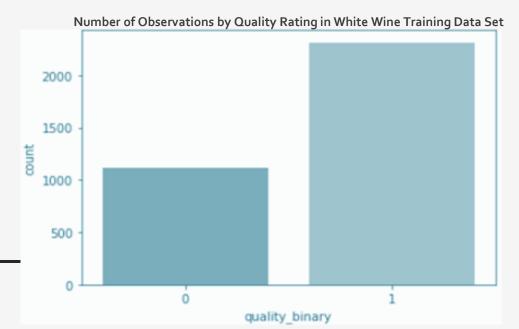
BINARY LOGISTIC REGRESSION MODEL WHITE WINE

1. Reclassified wine quality based on their quality score into a binary class:

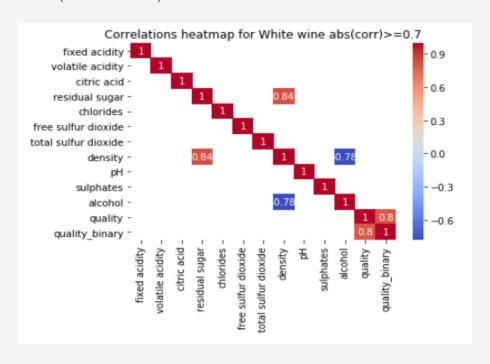


Results: more balance data.

Number of positive and negative responses are comparable



- 2. Checking Logistic Regression Conditions:
- Observations are independent of each other (given number of samples, assumed to be reasonable)
- There should be little or no multicollinearity among predictors, removed "Density"
- 3. Log transformed chlorides, free sulfur dioxide, and residual sugar and split data between Train and Test (ratio=0.3).



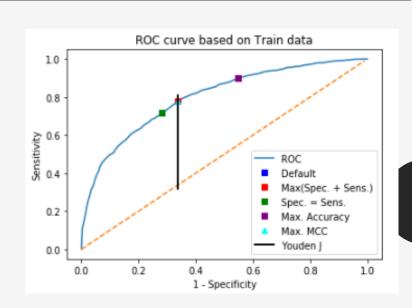
BINARY LOGISTIC REGRESSION MODEL WHITE WINE

Removed statistically insignificant predictors (citric acid and chlorides) after first iteration w/o much change to overall performance.

- The wine maker must balance risk of :
 - FP (predicted as premium but is economy) to avoid disappointing customers and loss of brand equity (higher precision)

 Predicted Value (y_hat)
 - FN (predicted as economy but is premium) to minimize loss of revenue/profit, may lead to higher volume/customers in the segment (higher recall).
 - Tune threshold value.
- Based on ROC curve, point of Max MCC was deemed to be optimum, with a threshold value of 0.62 replacing default value of 0.5 in the final model.
- Model performance with tune threshold value and Test data was "acceptable"

Data	Threshold	Precision	Recall	ACC	MCC
Train	0.50	0.774	0.898	0.754	0.400
Train	0.62	0.829	0.781	0.743	0.433
Test	0.62	0.796	0.775	0.729	0.416



731

187

True

Value

TP

FP

212

340

FN

TN

LINEAR REGRESSION, SVM

LINEAR REGRESSION

Independent variables:

- 1) fixed acidity
- 2) volatile acidity
- 3) citric acid
- 4) residual sugar
- 5) chlorides
- 6) free sulfur dioxide
- 7) total sulfur dioxide
- 8) density
- 9) pH
- 10) sulphates
- 11) alcohol

OLS Regression Resu	ılts		
Dep. Variable:	quality	R-equared:	0.295
Model:	OLS	Adj. R-equared:	0.293
Method:	Least Squares	F-statistic:	203.3
Date:	Wed, 20 Nov 2019	Prob (F-statistic):	0.00
Time:	07:30:52	Log-Likelihood:	-6464.3
No. Observations:	5847	AIC:	1.295e+04
Of Residuals:	5834	BIC:	1.304e+04
Of Model:	12		
Covariance Type:	nonrobust		

Dependant Variable:

Wine quality grouped as continuous scalar variable.

		uer atu err		E-Cld	[0.025	u.araj
cor	n et 109.76	34 16.374	6.704	0.000	77.665	141.862
fixed_acid	Ity 0.08	866 0.018	4.944	0.000	0.052	0.121
volatile_acid	Ity -1.48	867 0.085	-17.431	0.000	-1.654	-1.319
citric_ad	old -0.02	278 0.084	-0.333	0.740	-0.192	0.136
realdual_aug	gar 0.06	330 0.007	9.535	0.000	0.050	0.076
chlorid	98 -0.71	116 0.353	-2.013	0.044	-1.405	-0.019
free_sulfur_dlox	de 0.00	0.001	5.932	0.000	0.003	0.006
total_sulfur_dlox	de -0.00	0.000	-3.987	0.000	-0.002	-0.001
dena	Ity -108.87	08 16.595	-6.561	0.000	-141.403	-76.339
I	pH 0.50	0.098	5.139	0.000	0.311	0.695
aulphat	0.69	928 0.081	8.586	0.000	0.535	0.851
alcol	nol 0.21	31 0.021	10.274	0.000	0.172	0.254
wh	lte -0.36	889 0.062	-5.986	0.000	-0.490	-0.248
Omnibus:	130.269	Durbin-Wa	teon:	1.965		
Prob(Omnibus):	0.000 J	arque-Bera	(JB): 2	92.596		
Skew:	0.022	Prob	(JB): 2.9	91e-64		
Kurtoala:	4.095	Cond	No 32	6e+05		

Warnings:

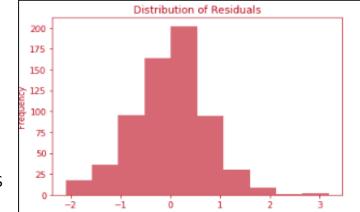
Highlights:

■ 90 / 10 – test train split

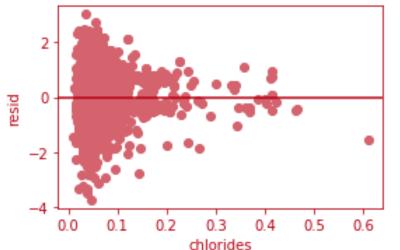
Overall score 0.3 R² on test

Observations:

- 1. weak Linearity amongst response and independent variables.
- 2. Independence of samples assumed as dataset is small.



Residuals versus chlorides



- 3. Residuals is normal
- 4. No constant variability in residuals plot.

^[1] Standard Errors assume that the covariance matrix of the errors is correctly specified

^[2] The condition number is large, 3.26e+05. This might indicate that there are strong multicollinearity or other numerical problems.

MODEL IMPROVEMENTS

Fine Tuning by backward elimination & p-value

Models have R² ranging from 0.26 to 0.3 which is not better than the first model.

Results of backward elimination

A				
model with removed feature	rsquared on test	aic	bic	fvalue
fixed_acidity volatile_acidity citric_acid residual sugar	0.3 0.27 0.31 0.29	12977.12 13249.51 12952.78 13043.08	13057.2 13329.6 13032.87 13123.17	218.65 184.55 221.77
chlorides free_sulfur_dioxide total sulfur dioxide	0.31 0.3 0.3	12956.73 12987.83 12968.58	13036.81 13067.91 13048.66	221.27 217.28 219.74
density pH sulphates	0.3 0.31 0.29	12995.65 12979.08 13026.09	13075.73 13059.16 13106.17	216.28 216.28 218.4 212.4
alcohol white	0.26 0.3	13057.51 12988.47	13137.59 13068.56	208.42

model with removed features	rsquared	aic	bic	fvalue
['citric_acid', 'chlorides']				

METHODS	RIDGE REGRESSION	NORMALIZATION BY LOG TRANSFORMATION	SVR WITH NORMALIZED DATA
	Used Regularization that constrains the magnitude of coefficients to reduce model complexity	 Free_sulfur_dioxide variable was transformed using np.sqrt. Chlorides, volatile_acidity and sulphates variables are log-transformed 	 Grid Search Cross Validation to find optimal SVM Params: Kernel, C, Gamma Data Normalization using Standard- Scaler
Outcome	Model coefficients are shrunk towards o	Statistically significant coefficients as almost all of the coefficients have p-value<0.05 p-values of citric_acid and chlorides are much better.	-
R ² on Test	0.31	0.33	0.45

CONCLUSIONS

- - White wine has more sugar
- INGREDIENTS THAT CAN SPOIL YOUR EXPERIENCE:
 - Volatile Acid turns to Vinegar
 - SO2 rotten egg like smell





RED OR WHITE?



AND THE WINNER IS...

(based on dataset used)

- SWEET FLAVOR
- FLORAL AROMA
- HIGHER ALCOHOL CONTENT

Best choice to make:

PORTUGUESE MADEIRA

SPANISH SHERRY

