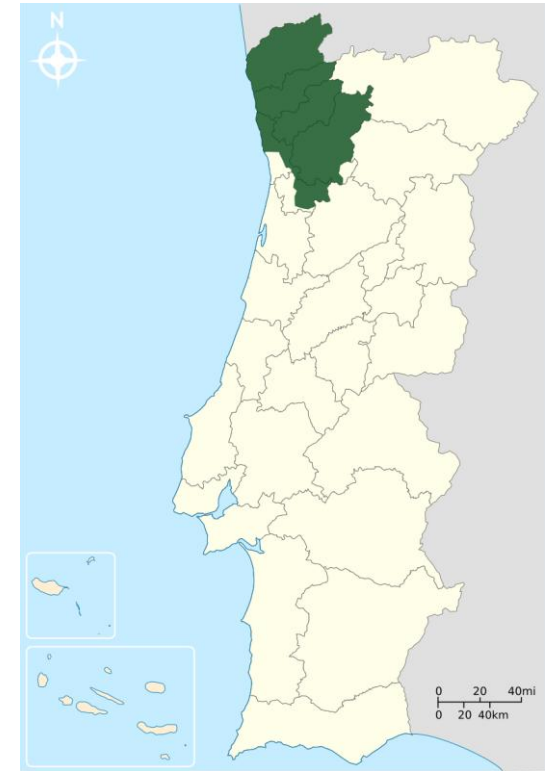# Clearing of the Vinho Verde Data Set



1. Duplicates
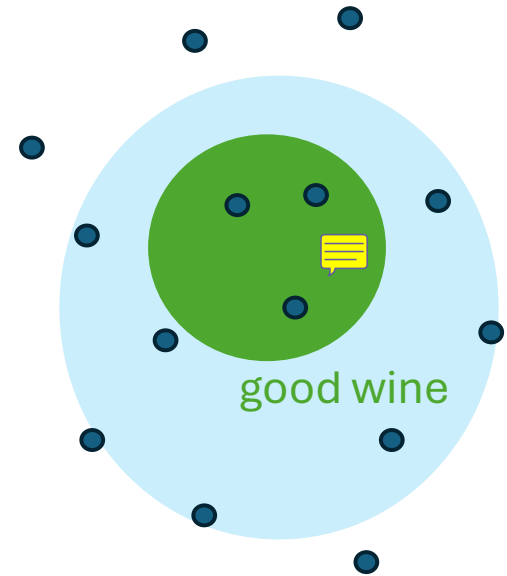2. Legal Limits
3. Outliers
4. Resulting Correlations

# Data Cleaning

- The data set is very noisy
  - Naturally: wine and taste is very complex and interrelated
  - Shallow

- Duplicates
  - Wines with exactly the same values
- „Illegal" Wines in EU
  - Total Sulfur Dioxide     < 210        mg/l
  - Volatile Acidity                  < 1.1        g/l
- Outliers Input Quantities
  -  values outside 1.5*IQR

Remove the above wines from the data set because
- Noisy and complex data set
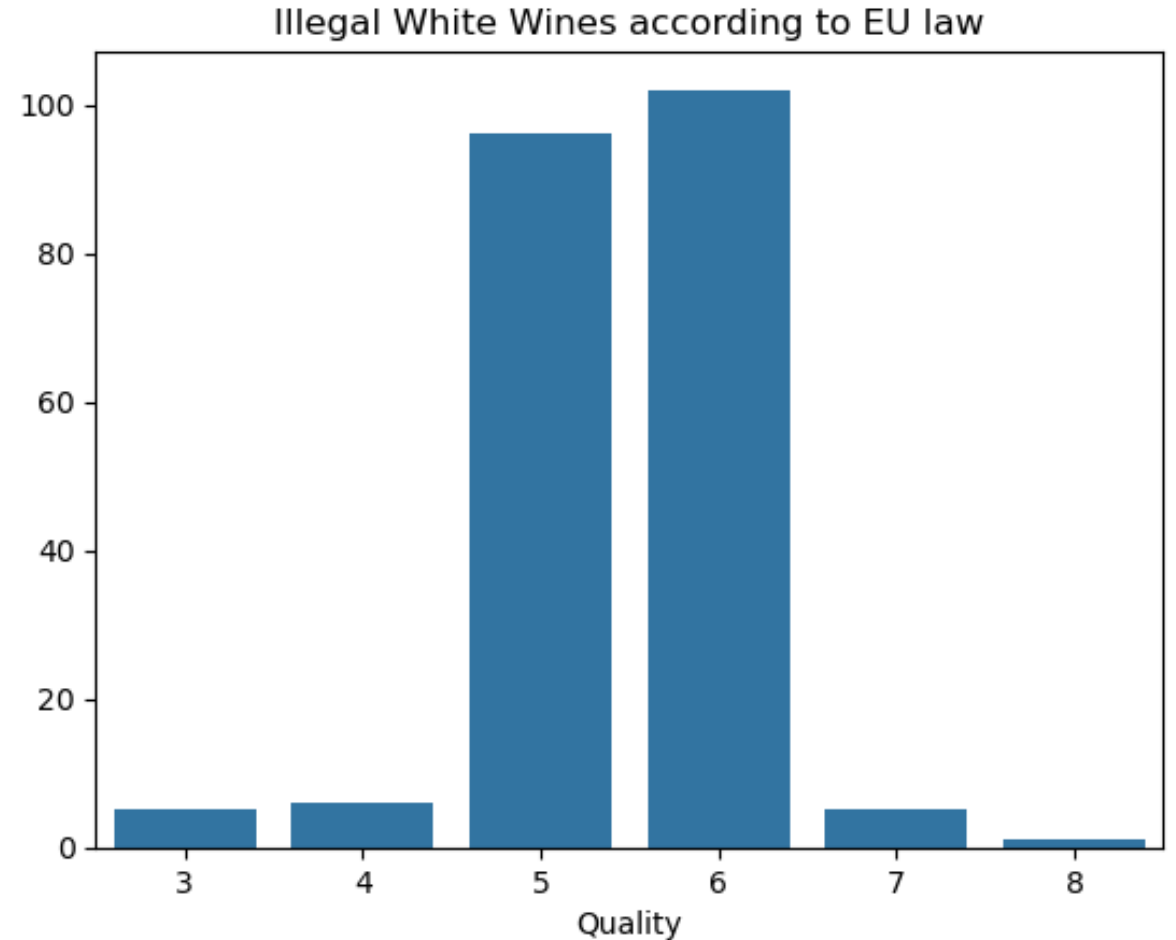- At least work on the main features

good wine

# Duplicates

- 937 of 4898 (19%) entries were dropped.
- Do not make sense.
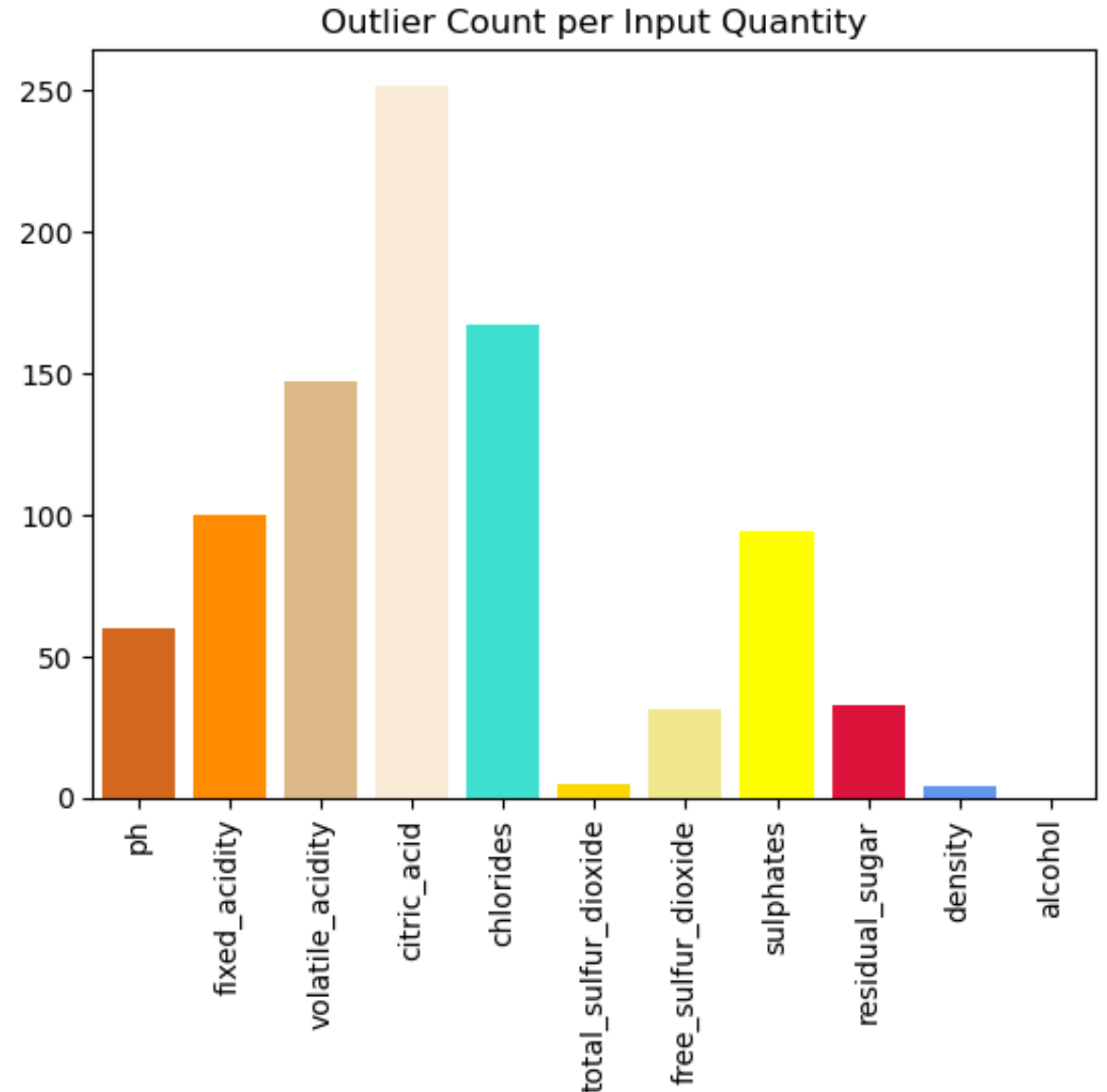- Effect: Duplicates shift correlation coefficients.

# Illegal Stuff

| Quantity | Limit | Count |
|---|---|---|
| Total Sulfur Dioxide | < 210mg/l | 215 |
| Volatile Acidity | < 1.1g/l | 0 |

- Remove if one of the legal limits is surpassed.
- Irrelevant, because we cannot sell those wines!



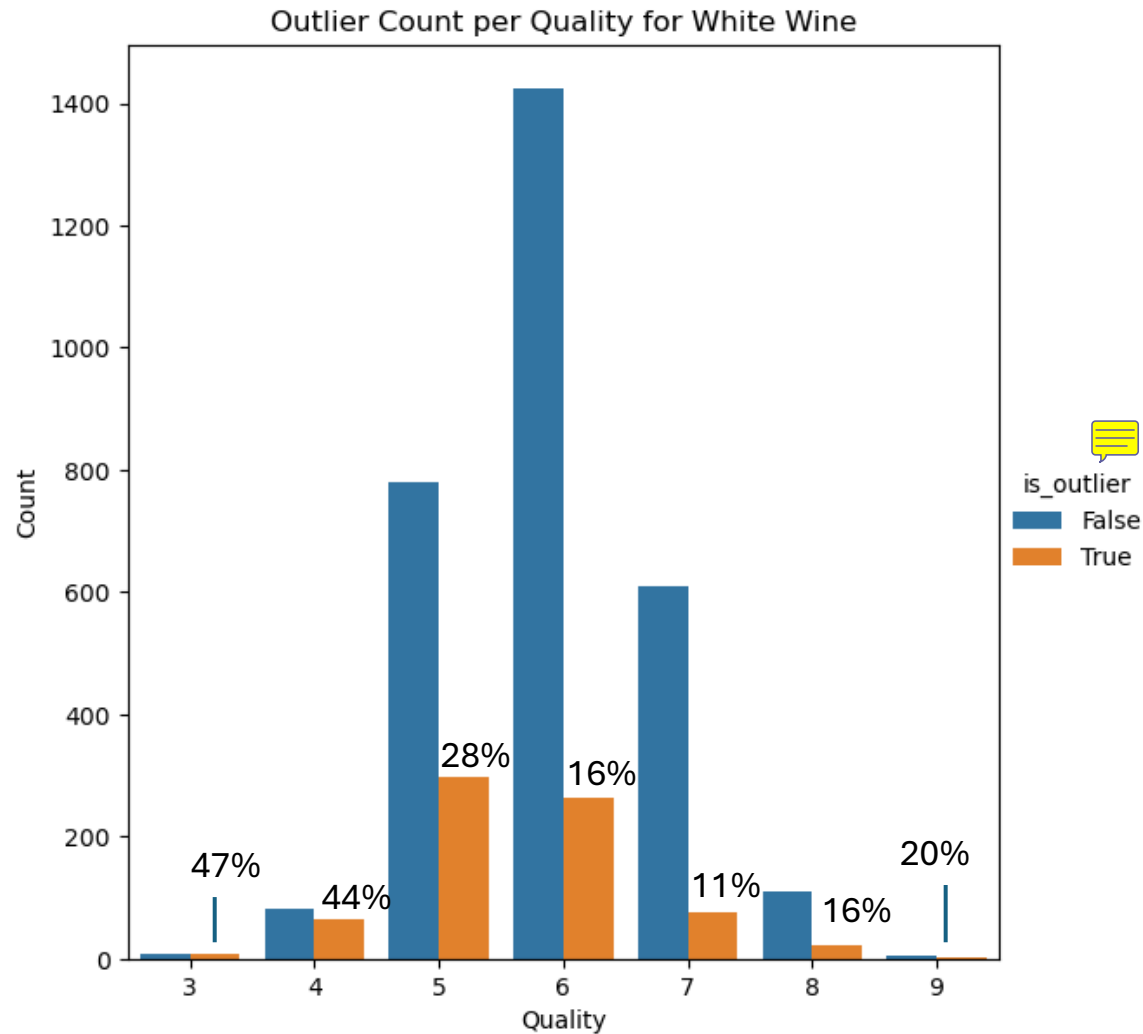Illegal White Wines according to EU law

# Outlier Wines

- 893 outliers in total

- **Acidity** seems to be a wi(l)dely varying quantity in wines.

- **Chlorides** (Terroir) varies strongly as well.

- **Sulfur Dioxide, Sulphates:** The use of preservatives varies a lot.



Outlier Count per Input Quantity

# Outlier Wines



Outlier Count per Quality for White Wine
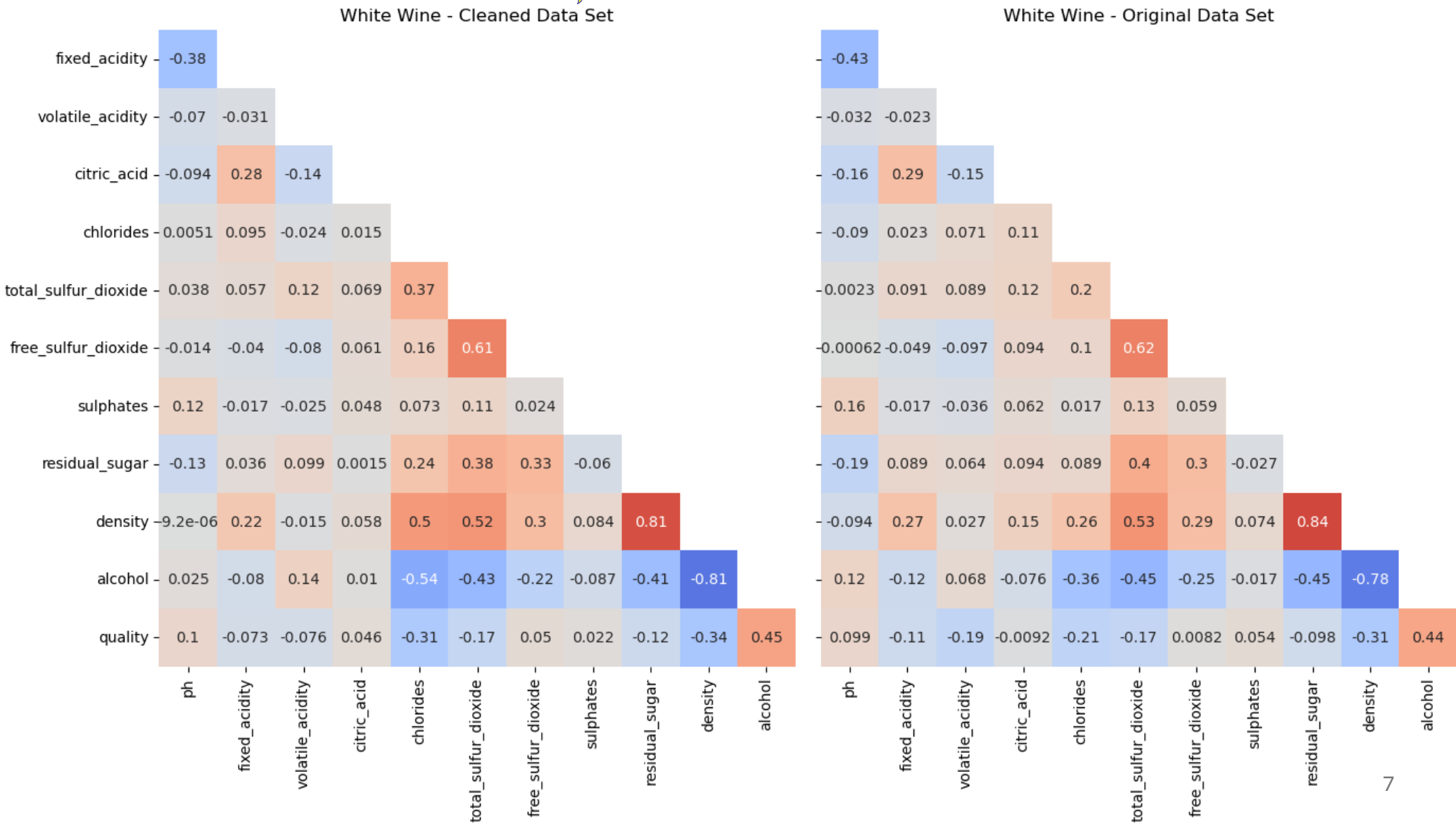
- Outlier Wine is a wine with at least one outlier in the input values

- Wine with
  - 1 outlier:          729
  - 2 outliers:        145
  - 3 outliers:        13

- First approach: Drop all outlier wines!

# Result



Correlations

White Wine - Cleaned Data Set

White Wine - Original Data Set

# Conclusions

- Data set with reduced noise mainly in Acidity
    - Acidity of wine is a very complex quantity and is not fully available in the data set.
    - It's good to get rid of it in order to find the main features of a good wine.

- Main features
    1. Alcohol
    2. Density
    3. Chlorides
    4. pH
    5. Residual Sugar

- Improvements
    - keep 5% of outliers per input quantity
    - keep the outliers in certain quantities, e.g. chloride
    - Select outliers more carefully (e.g.: take relation of sugar and density into account)
    - Use a different method to find outliers