# Summary

## Data exploratory

| Interest of the project

| Data wrangling / cleaning

| Exploratory data analysis (EDA)

## Machine learning

| Purpose

| Features selection
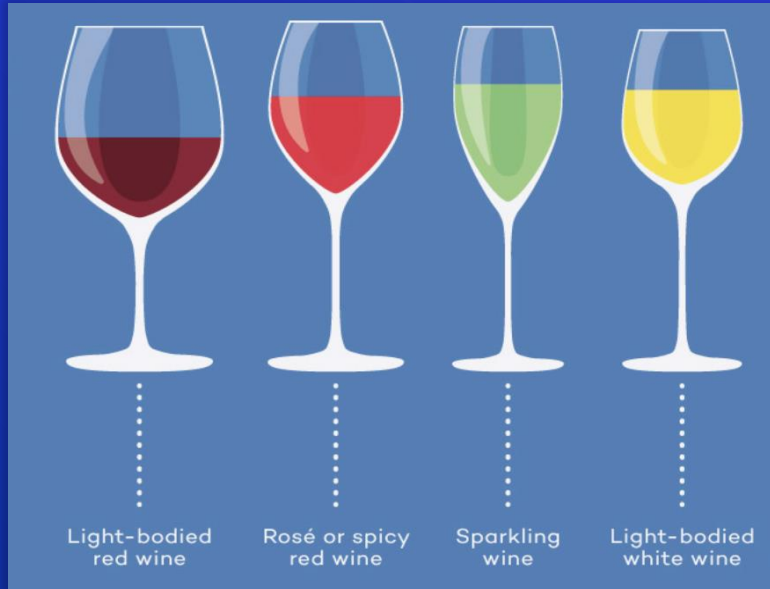
| Test set split for prediction

| Model selection

# Data exploratory

| Interest of the project

| Data wrangling / cleaning

| Exploratory data analysis (EDA)

# Interest of the project

Predict wine rating from features



40 000 🌐

3 € to 3000€

# Interest of the project

*The Vivino dataset*

San Francisco

**Features**

| Name of the bottle
| Country of provenance
| Region in the country
| Winery in the region
| Rating of the bottle 0 to 5 step 0.1
| Number of Ratings : the number of people
   which give a grade to the bottle
| Price of 1 bottle
| Year of production
| Type of the wine (Red, White, Rosé, Sparkling)

# Data wrangling / cleaning

Column created ✅

Encode variables ✅

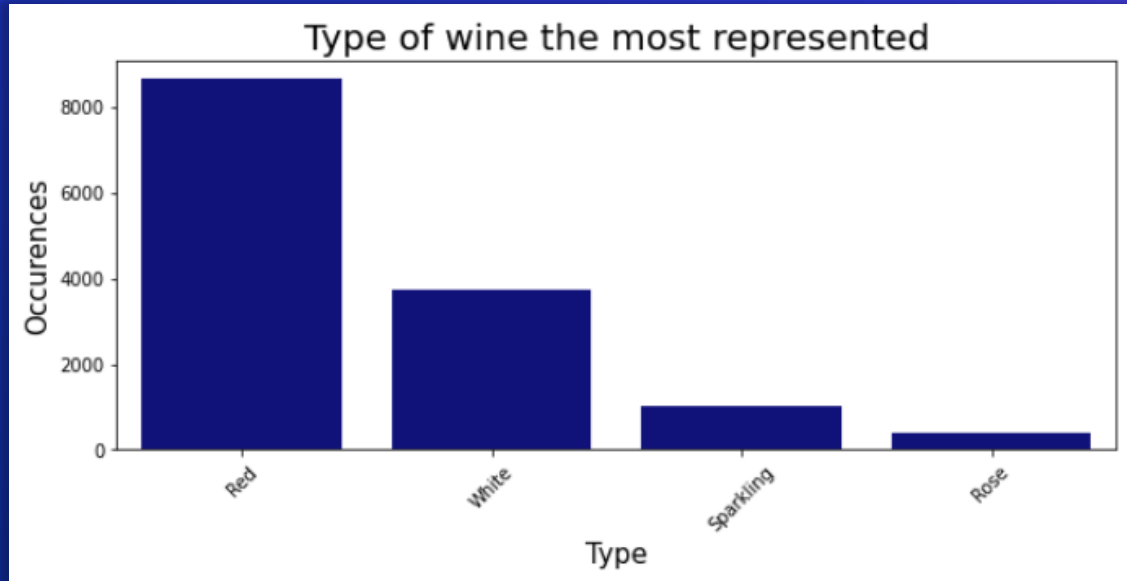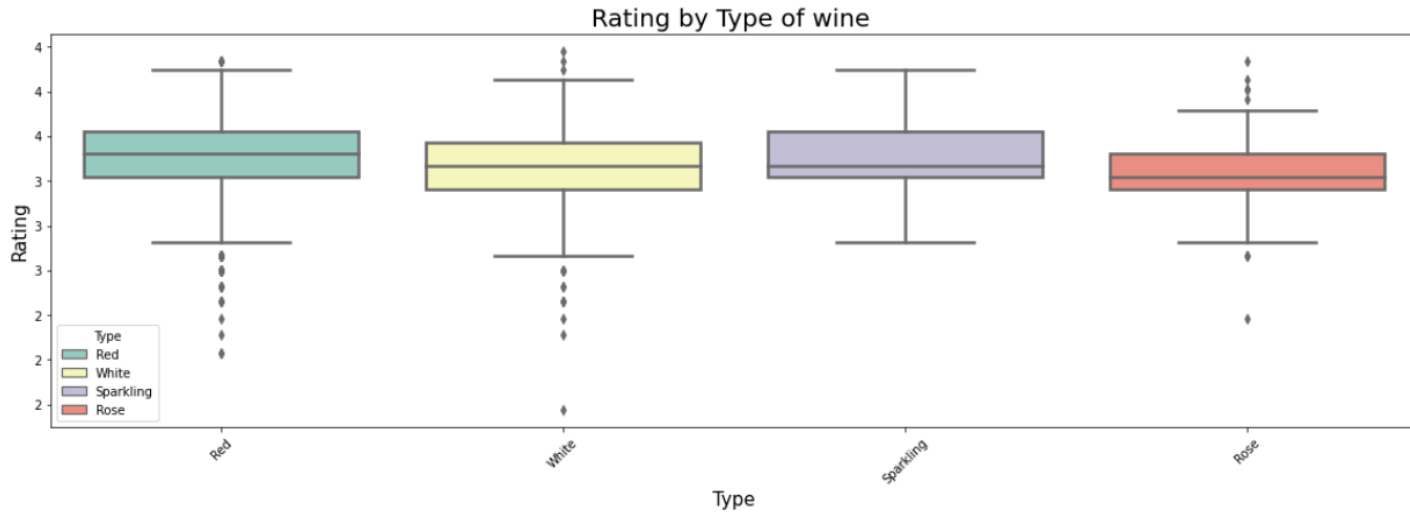**Features**
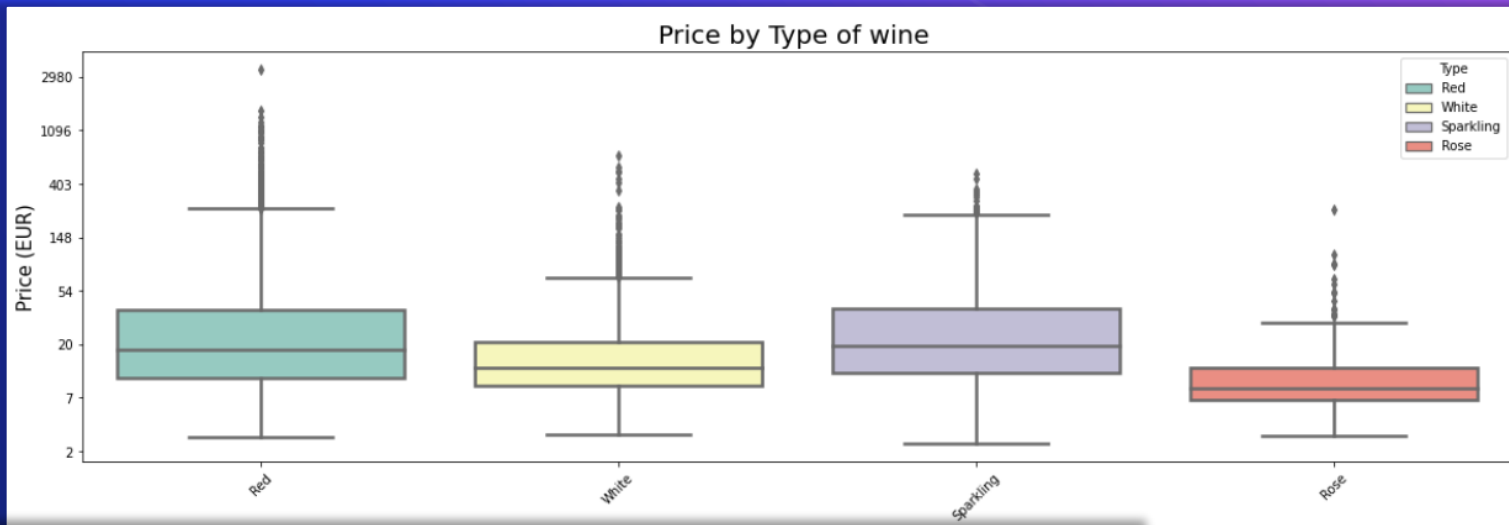
| Name
| Country
| Region
| Winery
| Rating
| NumberOfRatings
| Year
| Type

# Exploratory data analysis 📊



Type of wine the most represented

13 834
rows

# Boxplot



Price by Type of wine



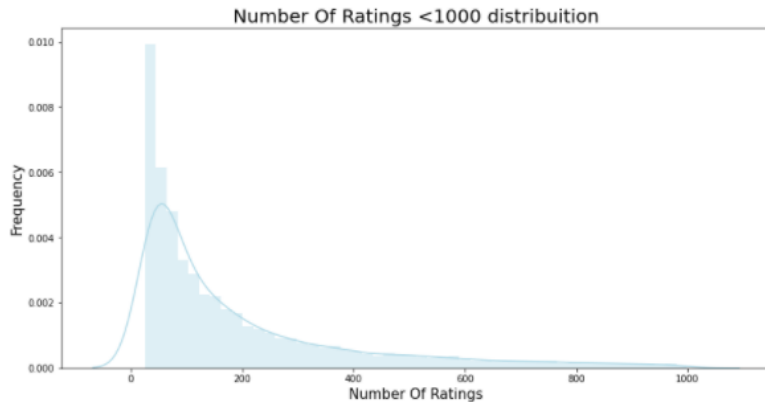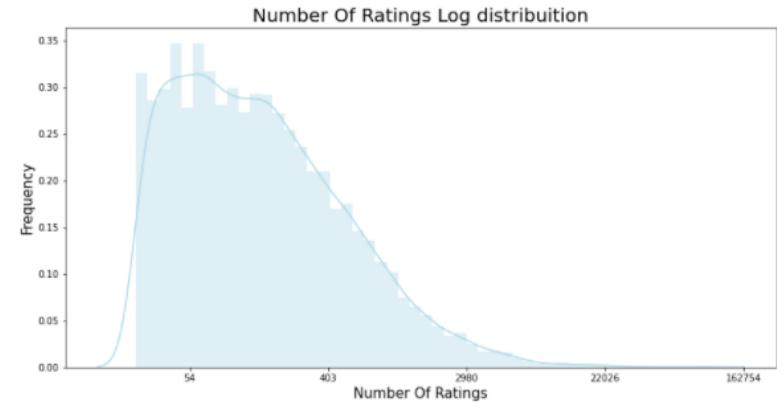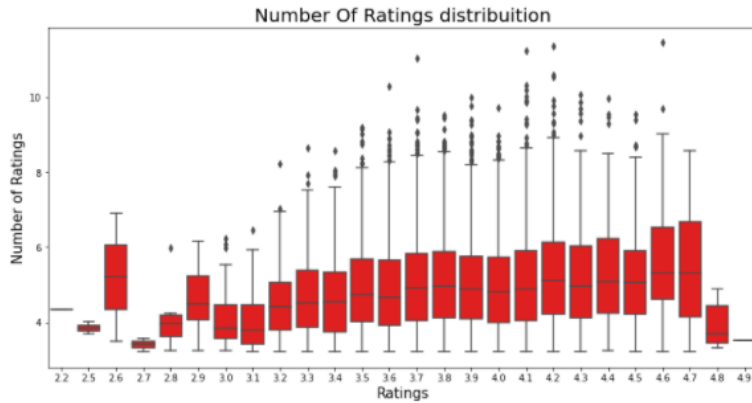Rating by Type of wine

8

# Correlation and trend





Rating by Price Distribution

# Rating plot

# Visualal Conclusion

"

| Number of Ratings has an exponential distribution

| For a large number of wines existing in Vivino, there is no Rating at all

| Problem for business ⚠️

# Machine learning

| Purpose

| Features engineering

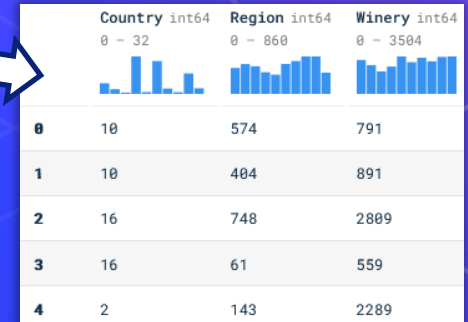| Test set split for prediction

| Model selection

# Purpose

1 Lack of Rating

2 Predict Rating

3 Inform the customer

# Feature engineering

**Features**

| Name
| Country (String)
| Region (String) ⟶ **Label Encoder**
| Winery (String)
| NumberOfRatings
| Price
| Year
| Type (Red, White, Sparkling, Rosé)
⟶ **One hot encoder**

# Test set split for prediction

- Definition of 3 set of data based on the **NumberOfRatings**
- We want to compare the performance between these sets

Number Of Ratings **LOW** : < 40

Number Of Ratings **MIDDLE** : > 40 & < 850

Number Of Ratings **HIGH** : > 850

15

# Model selection



scikit-learn algorithm cheat-sheet

# Comparison between Ridge Regressor and Ensembling model

## Ridge Regressor 😕

### Accuracy (MAE)

Number Of Ratings **LOW** : 0,22

⚠️ Number Of Ratings **MIDDLE** : 0,20

Number Of Ratings **HIGH** : 0,20

## Random Forest 😃

### Accuracy (MAE)

Number Of Ratings **LOW** : 0,17

Number Of Ratings **MIDDLE** : 0,13

Number Of Ratings **HIGH** : 0,11

# Thanks!

Any questions?