

# The Wranglers (Team28)

- (1) Naveen Kumar Conjeevaram Baskaran (email id: nc42@illinois.edu)
- (2) Alvin Do (email id: alvindo2@illinois.edu)
- (3) Hemantika Dasgupta (email id: hd8@illinois.edu)

## 1. Dataset Chosen

We opted for the Winery-Kaggle dataset (winemag-data-130k-v2.csv).

## 2. Description of Dataset

The winery-Kaggle dataset D contains about 130K records of wine reviews from around the world. This data was scraped from the Wine Enthusiast (<https://www.wineenthusiast.com/> - a popular site dedicated to wine, spirits, and the culture surrounding them) website during the week of June 15th, 2017.

The csv file originally consisted of 14 columns as mentioned below:

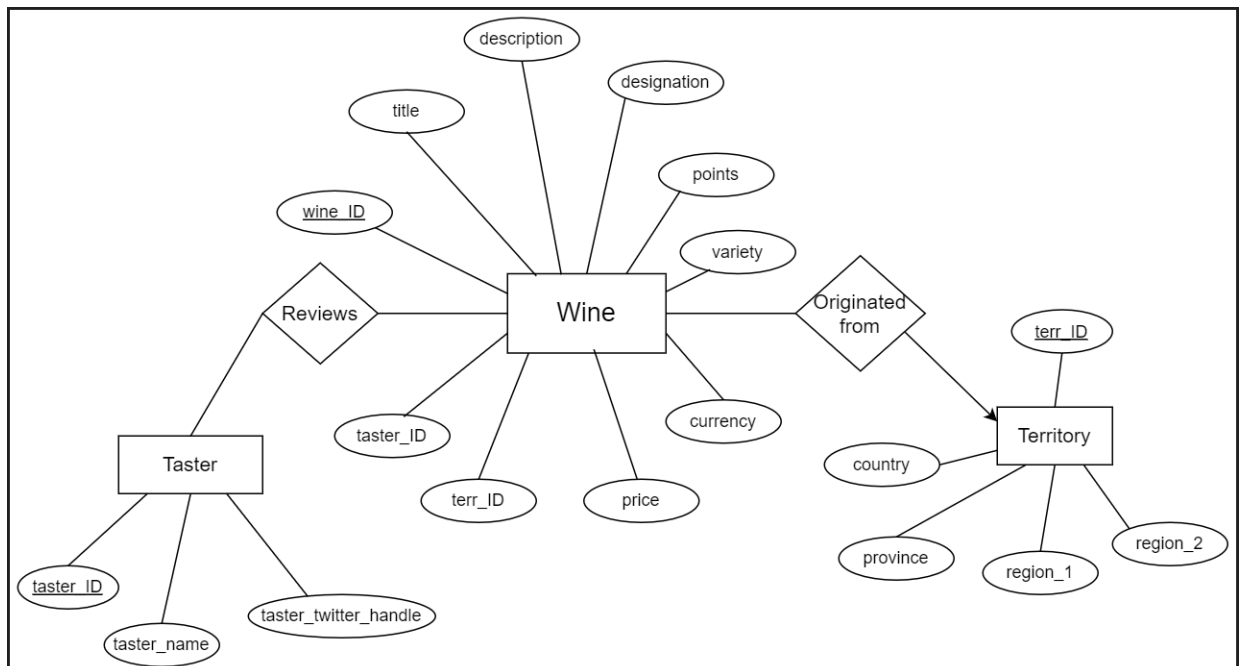
- i. *Wine\_ID*: A unique ID for each wine review record in the csv dataset.
- ii. *Country*: the country from which the wine was originated (ex: US).
- iii. *Description*: text describing the specific wine's taste, smell, look, feel, color, specialty, etc.
- iv. *Designation*: the vineyard within the winery where the grapes that made the wine are from.
- v. *Points*: number of points provided for the wine on a scale of 1-100 by the tasters tasting the wine.
- vi. *Price*: the cost for one bottle of the wine in local currency (ex: USD).
- vii. *Province*: the province or state within the country that originated the wine (ex: Washington).
- viii. *Region 1*: the area within the province or state (ex: Snipes Mountain).
- ix. *Region 2*: sometimes there are more specific regions specified within a wine growing area (ex:Columbia Valley inside of Snipes Mountain), but this value can sometimes be blank.
- x. *Taster Name*: name of the person who tasted and reviewed the wine.
- xi. *Taster Twitter Handle*: Twitter handle identifier for the person who tasted and reviewed the wine.
- xii. *Title*: name of the wine along with the territory which could be country, province or region if the province contains text "other".
- xiii. *Variety*: the type of grapes used to make the wine (ie Pinot Noir).
- xiv. *Winery*: the winery that produced the wine.

Additionally, some extra fields(as shown below) were added by us to provide more context around the data(currency) & to efficiently create a relationship model for this dataset(with additional IDs).

- xv. *Currency: The currency that's used to quote the wine price*
- xvi. *Taster\_ID: A unique identifier for the taster*
- xvii. *Terr\_ID: A unique identifier for the territory*

The dataset can be illustrated using the below ER diagram. The ER diagram can below be effectively transformed into relational tables, presenting a coherent and logical structure. There are 3 entities in the diagram below – Taster, Wine and Territory connected by the relationships – ‘Reviews’ and ‘Originated from’. As we could see from the data, multiple tasters can review multiple wines and provides the ratings. Therefore, we have a many-to-many relationship between the ‘Taster’ & ‘Wine’ entities. Many different wines can originate from a specific territory, hence we have a many-to-one relationship between the wine and territory entities.

We noticed that there's some information about the tasters and a lot of details of the wine in this dataset. Hence, we can organize the taster details in a separate ‘Taster’ table having the associated ID, name & twitter handle. The territory details of the wine specifically outlining the country, province, region can be moved to a separate ‘Territory’ table too. Lastly, all the wine attributes (title, description, designation, points, price, currency, variety etc.) can be put in a table of its own. The primary keys taster\_ID and terr\_ID from the Taster & Territory entities respectively are included in the ‘Wine’ entity as foreign keys for easier association of relational data.



### 3. Use Cases

a. “Main” use case U1: data cleaning is necessary and sufficient

Our target use case is that of a ‘wedding wine selection where the wedding is being held in the US and the guests would prefer European wine’. As we know, weddings are a grand affair and drinks play a crucial role. Hence, in order to arrange an assortment of wines that should delight the guests and yet be within the wedding budget, we need to look at the points, price, title, description, winery and territory details(country, province, region\_1, region\_2). Barring the ‘points’ data, every other data point mentioned above needs to be cleaned to some extent to make it ‘fit-for-purpose’ for the wedding wine analysis.

b. "Zero cleaning" use case U0: data cleaning is not necessary

If we were to just give a cursory look at this dataset to get a list of the wineries where they produce highly-rated wines and the wine variety, then just by looking at the 'points', 'winery' and 'variety' columns, we can draw an easy conclusion. The data in these 3 columns is pretty comprehensive and clean. There are no null values in 'points' and 'winery' columns and the 'variety' field has only 1 NULL value(out of the 130K records and can be easily researched), hence the data is good enough to use as is.

c. "Never enough" use case U2 : data cleaning is not sufficient

If we want to contact the tasters online using the information in this dataset, it would never fulfill this need as we noticed that a lot of records are missing the taster names and their twitter handles. No amount of manipulation in these 2 fields can point us to the correct individual, hence the data cleaning on these 2 columns will never be enough.

## 4. Data Quality Problems

One look at the dataset and we noticed the syntactic errors of various special characters and symbols that has corrupted a lot of data. Almost all text columns like description, province, region\_1, region\_2, title etc. are riddled with these symbols. Even the accented characters in some names(e.g. Saint-est  phe changed to Saint-Est  phe) have taken a different form in the csv data. Some data in the dataset also have the syntactic errors of white spaces. For e.g., an entry “Nicosia 2013 Vulk   Bianco (Etna)” in title column has 2 consecutive white spaces before the parenthesis. Moreover, we noticed a lot of semantic errors as well. There are a lot of missing data in columns like price, country and the regions. In our DB schema, we would define the price field to not have NULL values, hence few records in the price field would fail this NULL constraint. Additionally, we also noticed duplicate records in this dataset.

[illegible]

➔ Price field having NULL values, extra white spaces and junk characters in text fields

	country	description	points	price	province	region_1	region_2	title
1								
2	0 Italy	Aromas include tropical fruit, broo	87		Sicily & Sardin	Etna		Nicosia 2013 Vulk Bianco (Etna)
15	13 Italy	This is dominated by oak and oak-	87		Sicily & Sardin	Etna		Masseria Setteporte 2012 Rosso (Etna)
32	30 France	Red cherry fruit comes laced with l	86		Beaujolais	Beaujolais-Villages		Domaine de la Madone 2012 Nouveau (Beaujolais-Villages)
33	31 Italy	Merlot and Nero d'Avola form the	86		Sicily & Sardin	Sicilia		Duca di Salaparuta 2010 Calan Nero d'Avola-Merlot Red (Sicilia)
34	32 Italy	Part of the extended Calan Nero d'Avola	86		Sicily & Sardin	Sicilia		Duca di Salaparuta 2011 Calan Nero d'Avola-Merlot White (Sicilia)
52	50 Italy	This blend of Nero d'Avola and Syr	86		Sicily & Sardin	Sicilia		Viticultori Associati Canicatti 2008 Scialo Red (Sicilia)
56	54 Italy	A blend of Nero d'Avola and Nerell	85		Sicily & Sardin	Sicilia		Corvo 2010 Rosso Red (Sicilia)
81	79 Portugal	Grown on the sandy soil of Tejo, th	86		Tejo			Adega Cooperativa do Cartaxo 2014 Breda Touriga Nacional (Tejo)
139	137 South Afri	This is great Chenin Blanc, wood fe	90		Walker Bay			Beaumont 2005 Hope Marguerite Chenin Blanc (Walker Bay)
161	159 Italy	Intense aromas of ripe red berry, r	91		Tuscany	Brunello di Montalcino		Castello Romitorio 2011 Filo di Seta (Brunello di Montalcino)
165	163 France	Produced from vineyards donated	91		Beaujolais	Moulin-à-Vent		Collin-Bourisset 2011 Hospices Civils de Romanche (Moulin-à-Vent)
184	182 Italy	Loaded with bold, ripe fruit, exotic	88		Tuscany	Brunello di Montalcino		Abbadia Ardenga 2003 M. Vigna (Brunello di Montalcino)
196	194 Italy	Here's a traditional Chianti Classico	87		Tuscany	Chianti Classico		Campomaggio 2005 Chianti Classico
202	200 Italy	Aromas of mature black-skinned b	90		Tuscany	Brunello di Montalcino		Tenuta di Sesta 2011 Riserva (Brunello di Montalcino)
224	222 Italy	Chopped herb, forest floor, leather	90		Tuscany	Brunello di Montalcino		Talenti 2011 Trentennale (Brunello di Montalcino)
225	223 Italy	Bright and creamy, this savory whi	90		Northeastern	Alto Adige		Terlan 2015 Pinot Bianco (Alto Adige)
287	285 Austria	This is a very aromatic wine that's i	92		Kremstal			Salomon-Undhof 2011 Steiner Kellerberg Erste Lage Riesling (Kremstal)
290	288 Austria	While it feels rich and round, this i	92		Kremstal			Josef Schmid 2011 Kremser Gebirg Erste Lage Gruner Veltliner (Kremstal)
292	290 France	This is a wine that has great potent	92		Bordeaux	Saint-Estèphe		Château Lafon-Rochet 2011 Saint-Estèphe
293	291 Italy	This powerful Sagrantino opens wi	92		Central Italy	Sagrantino di Montefalco		Dionizio 2006 Sagrantino di Montefalco

## 5. Initial Plan for Phase-II

➔ The team discussed these data issues and came up with a draft plan as stated in the below steps:

- **Step 1(S1):** We'll use the winery-Kaggle dataset D that has a good collection of wine data and reviews to use for our target use case U1 of 'wedding wine selection where the wedding is being held in the US and the guests would prefer European wine'.
- **Step 2(S2):** We require the points, price, title, winery, description and territory details(country, province, region\_1, region\_2) from the dataset to get an idea of good-rated wine within a budget and a certain country we would like it from. We don't require the taster name and the twitter handle for this analysis. Alongside from the visual inspection & analysis we performed on the dataset, we are also planning to use 'Regex' to identify the syntactic errors and erroneous data patterns and 'Datalog' to identify the integrity constraint violations.
- **Step 3(S3):** To help with this arduous task of data cleaning, we are planning to use (i) **Regex** to identify the syntactic errors (special symbols, characters, white spaces), (ii) **OpenRefine** to clean the syntactic errors to transform the dataset. Since OpenRefine allows use of GREL, Regex and Python languages, we intend to leverage these options to perform our data cleaning. (iii) **Datalog** to check for integrity constraints like duplicate data & NULL values (iv) **SQLite** to fix semantic errors and update the records with values that can be determined. As the weeks goes by and we learn about other tools such as YesWorkflow, we would also explore possibilities to use this in our data cleaning exercise.

- **Step 4(S4):** Once the new dataset is ready, we want to compare this with the original dataset. We plan to use SQLite queries to do this 'before & after' comparison. For e.g., we plan to write queries for checking counts of NULL values and verify if duplicates have been removed from the dataset.
  - **Step 5(S5):** As we understand, documentation of the steps will greatly help with this assignment evaluation. Hence, we are thinking of using the OpenRefine JSON log to capture the cleaning steps that will occur in OpenRefine tool and for the rest of the items, we can prepare a document summarizing the steps and capturing screenshots of the results.
- ➔ Steps S1, S2, S4 and S5 will be a joint effort as we have a weekly meeting cadence and we would address the work during that time. We plan to dedicate 2-3 hrs per week till the week of July 24 and finish the project.
- ➔ For the tasks listed in step S3, Alvin will be responsible for the Regex and Datalog work; Naveen will be responsible for the OpenRefine tool usage and Hemantika will own the SQLite work.