

MS Azure ML Designer를 활용한 개인 실습

Lab : 와인 데이터

강명호

배경 및 모델링 목표

배경 및 목표

Wine Quality 데이터 세트를 분석하여 와인의 종류(white or red)를 분류하고,
와인의 특성을 통해 품질을 예측하는 모델 구현

- Story :

- 와인 메이커는 와인 분석가를 고용하여 산도, 단맛, 알코올 수준 등 생산 중인 와인의 특성을 측정하고 개선함
- 최근, 와인 분석가가 다른 프로젝트에 배정되어서 와이너리의 와인 분석 작업을 진행하지 못하게 되었음
- 와인 분석가를 대체하기 위한 머신러닝 모델이 필요한 상황임



배경 및 목표

Wine Quality 데이터 세트를 분석하여 와인의 종류(white or red)를 분류하고,
와인의 특성을 통해 품질을 예측하는 모델 구현

- 데이터 소스 : UCI Machine Learning Repository
- 포르투갈의 Vinho Verde의
레드 및 화이트 와인 데이터에서
일부 발췌 (민감한 정보 제외)



와인 데이터 세트를 이용하여 모델 구현

- 와인의 종류를 분류하는 모델, 또는
- 와인의 품질을 측정하는 모델



데이터 수집

데이터 수집 (참고)



Wine Quality

Donated on 10/6/2009

Two datasets are included, related to red and white vinho verde wine samples, from the north of Portugal. The goal is to model wine quality based on physicochemical tests (see [Cortez et al., 2009],...)

Dataset Characteristics Multivariate	Subject Area Business	Associated Tasks Classification, Regression	1 citations 369950 views
Feature Type Real	# Instances 4898	# Features 11	Keywords Chemistry

Dataset Information

Additional Information

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine. For more details, consult: <http://www.vinhoverde.pt/en/> or the reference [Cortez et al., 2009]. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

Creators

- Paulo Cortez
- A. Cerdeira
- F. Almeida
- T. Matos
- J. Reis

These datasets can be viewed as classification or regression tasks. The classes are ordered and not balanced (e.g. there are many more normal wines than excellent or poor ones). Outlier detection

DOWNLOAD

IMPORT IN PYTHON

CITE

데이터 수집 (참고)

三

winequality.names

winequality-red.csv

winequality-white.csv

데이터 수집 (참고)

三

winequality.names

 winequality-red.csv

winequality-white.csv

데이터 수집

Kaggle에 접속하여 Wine Quality 데이터 세트 검색

The screenshot shows the Kaggle homepage. On the left is a sidebar with the user's profile name "kaggle" and a navigation menu including "Create", "Home", "Competitions", "Datasets", "Models", "Code", "Discussions", "Learn", and "More". Below this is a section for "Your Work" and a "VIEWED" section. The main content area features a search bar at the top right with the placeholder "Search". A large "Welcome, Myeongho Kang!" message is centered, followed by the text "You're on a roll! Jump back in, or start something new.". To the right, there is a "LOGIN STREAK" indicator showing "4 days a new record!". Below this are four summary cards: "Datasets" (0 total created), "Notebooks" (0 total created), "Competitions" (0 total joined), and "Discussions" (0 total posted). At the bottom, there is a call-to-action "How to start: Choose a focus for today" with the sub-instruction "Help us make relevant suggestions for you".

데이터 수집

Kaggle에 접속하여 Wine Quality 데이터 세트 검색

← Wine Quality

Notebooks 3,938 Comments 2,540 Datasets 242 Topics 230 Competitions 77 Models 3

Filter by 7,030 Results

DATE

- Last 90 days 267
- This week 12
- Today 1

VIEWED BY YOU

- Viewed 2
- Not Viewed 7,028

 Red Wine Quality
Notebook · 7mo ago · by [Nima Pourmoradi](#)
Import data by using pandas liblary and using read_csv method data = pd.read_csv('/kaggle/input/re

 Red Wine Quality ~ EDA & Classification
Notebook · 2y ago · by [Mustanger](#)
>6 - 8</mark> ##### - Low quality wine: <mark>3 - 5</mark> df["quality"] = np.w

 Red Wine Quality Prediction
Notebook · 2y ago · by [Halime Doğan](#)
`/input/red-wine-quality-cortez-et-al-2009/winequality-red.csv") df.head() def grab_col_names(data`

데이터 수집

Kaggle에 접속하여 Wine Quality 데이터 세트 검색

← Wine Quality

← Wine Quality

Notebooks 3,938 Comments 2,540 Datasets 242 Topics 230 Competitions 77 Models 3

Filter by Relevance

DATE

- Last 90 days 17
- Viewed 4
- Not Viewed 238

VIEWED BY YOU

CREATOR

- You 0
- Others 242

242 Results

 **Red Wine Quality**
Dataset - 7y ago · by UCI Machine Learning
brand, wine selling price, etc.). 2806 245,337 downloads

 **Wine Quality Dataset**
Dataset - 2y ago · by M Yasser H
Wine Quality Prediction - Classification Prediction 641 55,603 downloads

 **Wine Quality**
Dataset - 6y ago · by Raj Parmar
brand, wine selling price, etc.). 259 41,175 downloads

 **Wine Quality** 38

A red arrow points to the third dataset entry, "Wine Quality" by Raj Parmar.

데이터 수집



RAJ PARMAR · UPDATED 6 YEARS AGO

▲ 259 New Notebook [Download \(100 kB\)](#) ▾

Wine Quality



Data Card Code (171) Discussion (3) Suggestions (0)

About Dataset

Data Set Information:

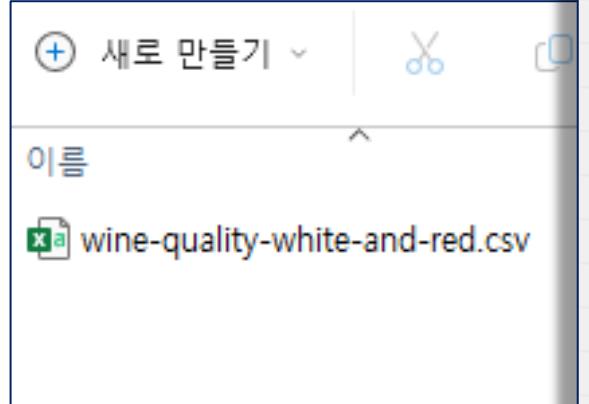
The dataset was downloaded from the UCI Machine Learning Repository.

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine. The reference [Cortez et al., 2009]. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

These datasets can be viewed as classification or regression tasks. The classes are ordered and not balanced (e.g. there are much more

Usability ⓘ	7.06
License	Other (specified in description)
Expected update frequency	Not specified
Tags	Classification, Regression, Wine, Physicochemical, Sensory, Portuguese, Vinho Verde, Red, White, Variants, UCI Machine Learning Repository, Cortez et al., 2009, Privacy, Logistic Issues, Inputs, Outputs, Classes, Ordered, Not Balanced, Data Types, Grape Types, Wine Brand, Selling Price, Classification Task, Regression Task

데이터 수집



About Dataset

Data Set Information:

The dataset was downloaded from the UCI Machine Learning Repository.

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine. The reference [Cortez et al., 2009]. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

These datasets can be viewed as classification or regression tasks. The classes are ordered and not balanced (e.g. there are much more normal wines than excellent or poor ones). Outlier detection algorithms could be used to detect the few excellent or poor wines. Also, we are not sure if all input variables are relevant. So it could be interesting to test feature selection methods.

Two datasets were combined and few values were randomly removed.

Attribute Information:

For more information, read [Cortez et al., 2009].

Usability i

7.06

License

Other (specified in description)

Expected update frequency

Not specified

Tags

Earth and Nature

Classification

Alcohol

Regression

데이터 수집

- fixed acidity : 고정 산도
- volatile acidity : 휘발성 산도
- citric acid : 시트르산
- residual sugar : 잔류 당분
- chlorides : 염화물
- free sulfur dioxide : 자유 이산화황
- total sulfur dioxide : 총 이산화황
- density : 밀도
- pH
- sulphates : 황산염
- alcohol
- quality : 0 ~ 10(높을 수록 좋은 품질)

Attribute Information:

For more information, read [Cortez et al., 2009].

Input variables (based on physicochemical tests):

- 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

Output variable (based on sensory data):

- 12 - quality (score between 0 and 10)