

Joonho (Joshua) Kim | U77524942

Prof. Alishah Chator

CDS DS210 B1

December 10, 2024

Wine Quality Prediction Using Rust

Project Overview

This project focuses on predicting the quality of wines (red and white) based on their chemical properties. Using the UC Irvine Machine Learning Repository's Wine Quality Dataset, the project analyzes the features of wines to build predictive models and generate insights into the factors that influence wine quality. The dataset contains chemical attributes such as alcohol content, pH, residual sugar, and quality scores on a scale of 0 to 10.

The project utilizes Rust to implement the solution, emphasizing efficiency, modularity, and reliability.

Dataset Details

Name:

UC Irvine Machine Learning Repository: "Wine Quality Dataset"

Source:

<https://archive.ics.uci.edu/dataset/186/wine+quality>

Size:

4898 entries

Why This Dataset Was Chosen:

1. Predictive Insights:

The dataset's chemical properties provide a strong basis for predicting wine quality using machine learning.

2. Comparative Analysis:

It contains data for both red and white wines, which enables an intriguing comparative study.

3. Manageable Size:

The dataset is large enough for meaningful analysis but remains computationally feasible in Rust.

Problem Statement

Objective:

To predict wine quality based on chemical attributes and identify key factors influencing the quality of red and white wines.

Key Questions:

1. What are the average quality and alcohol content of red and white wines?
2. Which chemical features most significantly affect wine quality?
3. How do the predictive models perform for red and white wines separately?

4. What differences exist in the importance of features for red and white wines?

Project Structure & Components

Modules:

1. Data Ingestion:
 - Parses CSV files and loads wine data into appropriate data structures.
 - Normalizes the features for consistent analysis.
2. Exploratory Data Analysis (EDA):
 - Calculates summary statistics like average alcohol content and quality.
3. Feature Importance:
 - Implement logistic regression to evaluate the importance of features for quality prediction.
4. Model Implementation:
 - Predicts wine quality using logistic regression and evaluates the predictions.
5. Comparison Module:
 - Compares feature importance and model performance across red and white wines.
6. Tests:
 - Verifies the correctness of key functionalities such as data parsing, normalization, and feature importance computation.

Milestones:

1. Week 1: Data ingestion and cleaning
2. Week 2: Implement exploratory data analysis (EDA)

3. Week 3: Build predictive models and compute feature importance
4. Week 4: Final testing, evaluation, and documentation

Code Analysis

1. Data Ingestion (data_ingestion.rs)

Code Overview:

The data_ingestion.rs module is responsible for loading, cleaning, and normalizing the dataset. It reads the data from the CSV files for red and white wines and processes each feature for later analysis.

Key Functions:

- load_data(file_path: &str):
 - Reads the CSV file into a vector of hashmaps where each key is a column name (e.g., alcohol, quality) and each value is the corresponding value in the row.
 - Ensures that the data is parsed correctly into string-keyed hashmaps for generality.
- normalize_features(data: &mut [HashMap<String, String>]):
 - Normalizes each feature by scaling values between 0 and 1, using the formula:
$$\text{normalized value} = (x - \text{min}) / (\text{max} - \text{min})$$
 - This ensures all features have the same scale, which is essential for algorithms like logistic regression.

Analysis:

- The normalization ensures comparability between features with different ranges (e.g., residual sugar vs pH).
- A potential improvement could include error handling for missing or invalid data during normalization.
- The modularity of this code makes it reusable for datasets with similar structures.

2. Exploratory Data Analysis (eda.rs)

Code Overview:

This module computes basic descriptive statistics for the dataset, such as average alcohol content and quality ratings for both red and white wines.

Key Functions:

- `compute_statistics(data: &[HashMap<String, String>])`:
 - Iterates through the dataset to compute averages for numerical columns like alcohol and quality.
 - Uses pattern matching to ensure only valid numerical fields are processed.

Analysis:

- Computing averages provides an initial understanding of the dataset and allows for meaningful comparisons (e.g., average alcohol content of red vs white wines).

- Additional EDA steps—such as plotting histograms or correlations between features—could further enrich the analysis.

3. Feature Importance (feature_importance.rs)

Code Overview:

The `feature_importance.rs` module calculates the importance of each feature in predicting wine quality. It uses a simplistic placeholder computation based on feature correlations with quality.

Key Functions:

- `compute_feature_importance(data: &[HashMap<String, String>], weights: &[f64]):`
 - The `weights` parameter represents the coefficients learned by the logistic regression model.
 - Maps each feature to its corresponding weight to rank their importance.

Analysis:

- The feature importance is directly tied to the predictive model. Since weights are derived from logistic regression, features with higher absolute weights contribute more to predictions.
- An observed issue is that the feature importance values are all zero.

4. Logistic Regression Model (model.rs)

Code Overview:

This module implements logistic regression for predicting wine quality.

Key Functions:

- `train_logistic_model(data: &[HashMap<String, String>])`:
 - Uses gradient descent to iteratively update weights based on the prediction error.
 - Applies the logistic function to map predictions to probabilities.
- `predict(data: &[HashMap<String, String>], weights: &[f64])`:
 - Computes the dot product of weights and input features to generate predictions.

Analysis:

- The implementation uses a straightforward approach to logistic regression; however, the gradient computation ($\text{error} * x$) may require refinement to handle edge cases or numerical instability.
- Including regularization (e.g., L2 regularization) could improve the model's generalization and reduce overfitting.

5. Feature Comparison (comparison.rs)

Code Overview:

This module compares the feature importance between red and white wines,

which identifies features that significantly differ in their contributions to quality predictions.

Key Functions:

- `compare_features(red_importance: &[f64], white_importance: &[f64])`:
 - Computes the difference between corresponding features for red and white wines.
 - Outputs a sorted list of features by their differences in importance.

Analysis:

- This comparison provides insights into how different chemical properties influence the quality of red vs. white wines. For example, alcohol content might matter more for red wines, whereas residual sugar could dominate in whites.
- The comparison logic is correct, but the zeroed-out importance values make it less informative.

6. Main Program (main.rs)

Code Overview:

The main.rs file orchestrates the overall workflow:

1. Loads and normalizes data.
2. Computes statistics and runs EDA.
3. Trains logistic regression models for red and white wines.
4. Calculates and compares feature importances.

Key Sections:

- Data Loading and Normalization:
Ensures that both red and white datasets are ready for analysis.
- EDA Results:
Outputs basic statistics, which provides an overview of the dataset.
- Model Training:
Trains separate logistic regression models for red and white wines, which generates weights for feature importance.
- Output:
Displays EDA results, feature importances, and predictions.

Analysis:

- The modular design improves readability and maintainability.
- Better error handling (e.g., for empty datasets or division by zero in normalization) could enhance robustness.

7. Tests (tests/mod.rs)

Code Overview:

The test module validates the functionality of key components.

Key Test Cases:

- Data Ingestion Tests:
Ensure correct parsing and normalization of input data.
- Logistic Regression Tests:
Check that weight updates occur as expected during gradient descent.
- Prediction Accuracy:
Verify predictions against a small, manually computed dataset.

Analysis:

- These tests cover core functionality but could be expanded to include edge cases (e.g., datasets with missing or invalid values).
- Benchmarking tests could measure the efficiency of normalization and gradient descent.

Output Analysis

Results:

1. Red Wine Statistics:
 - Average Alcohol: 0.3112 (normalized)
 - Average Quality: 0.5272 (normalized)
2. White Wine Statistics:
 - Average Alcohol: 0.4055 (normalized)
 - Average Quality: 0.4797 (normalized)
3. Feature Importance:

- Red Wine:

All feature importances were calculated as 0.0.

- White Wine:

Similarly, all features had the importance of 0.0.

4. Predictions: All predictions were 0.0.

RUST_BACKTRACE=full cargo run:

```
Red wine statistics: {"avg_quality": 0.527204502814258, "avg_alcohol": 0.3112281844903063}
White wine statistics: {"avg_alcohol": 0.40552689771335115, "avg_quality": 0.4796515183748503}
thread 'main' panicked at src/feature_importance.rs:7:53:
called `Option::unwrap()` on a `None` value
stack backtrace:
0:      0x1002e4b40 - std::backtrace_rs::backtrace::libunwind::trace::hbebc8679d47bdc2c
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/.../backtrace/src/backtrace/libunwind.rs:116:5
1:      0x1002e4b40 - std::backtrace_rs::backtrace::trace_unsynchronized::h3a2e9637943241aa
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/.../backtrace/src/backtrace/mod.rs:66:5
2:      0x1002e4b40 - std::sys::backtrace::_print_fmt::he430849680584674
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/sys/backtrace.rs:65:5
3:      0x1002e4b40 - <std::sys::backtrace::BacktraceLock::print::DisplayBacktrace as core::fmt::Display>::fmt::h243268f17d714c7f
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/sys/backtrace.rs:40:26
4:      0x1002fcc8c - core::fmt::rt::Argument::fmt::h0d339881c25f3c31
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/fmt/rt.rs:173:76
5:      0x1002fcc8c - core::fmt::write::hb3cfb8a30e72d7ff
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/fmt/mod.rs:1182:21
6:      0x1002e2f08 - std::io::Write::write_fmt::hfbb2314975de9ecf1
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/io/mod.rs:1827:15
7:      0x1002e5c1c - std::sys::backtrace::BacktraceLock::print::he14461129ccbfef5
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/sys/backtrace.rs:43:9
8:      0x1002e5c1c - std::panicking::default_hook::{{closure}}::h14c7718ccf39d316
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:269:22
9:      0x1002e5840 - std::panicking::default_hook::hc62e60da3be2f352
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:296:9
10:     0x1002e6630 - std::panicking::rust_panic_with_hook::h09e8a656f11e82b2
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:800:13
11:     0x1002e6008 - std::panicking::begin_panic_handler::{{closure}}::h1230eb3cc91b241c
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:667:13
12:     0x1002e4fcc - std::sys::backtrace::_rust_end_short_backtrace::hc3491307aceda2c2
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/sys/backtrace.rs:168:18
13:     0x1002e5cf8 - rust_begin_unwind
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:665:5
14:     0x10030446c - core::panicking::panic_fmt::ha4b80a05b9fff47a
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/panicking.rs:74:14
15:     0x1003044d8 - core::panicking::panic::h298549a7412a7069
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/panicking.rs:148:5
16:     0x1003043f4 - core::option::unwrap_failed::hb7af631ec4f78cd6
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/option.rs:2020:5
17:     0x1002bf864 - core::option::Option<T>::unwrap::h8098096df415fd0d
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/option.rs:970:21
18:     0x1002bf864 - wine_quality_project::feature_importance::compute_feature_importance::{{closure}}::h052358b2bb96c310
    at /Users/USER/DS210_final_project/wine_quality_project/src/feature_importance.rs:7:31
19:     0x1002c2cfc - alloc::slice::<impl1 [T]>::sort_by::{{closure}}::h750c29b35053491c
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/alloc/src/slice.rs:289:34
20:     0x1002b229c - core::slice::sort::shared::smallsort::insert_tail::h159fa46c03423119
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/slice/sort/shared/smallsort.rs:545:13
21:     0x1002b2f44 - core::slice::sort::shared::smallsort::insertion_sort_shift_left::hca29baa5d3ac704
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/slice/sort/shared/smallsort.rs:600:13
22:     0x1002c2c44 - core::slice::sort::stable::sort::h0d01738a577f65f2
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/slice/sort/stable/mod.rs:43:9
23:     0x1002c2c44 - alloc::slice::stable_sort::h3225f49f6c91b28e
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/alloc/src/slice.rs:883:5
24:     0x1002c2c94 - alloc::slice::<impl1 [T]>::sort_by::haa6da045214abe94
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/alloc/src/slice.rs:289:9
25:     0x1002a94dc - wine_quality_project::feature_importance::compute_feature_importance::haf03f11741e23803
    at /Users/USER/DS210_final_project/wine_quality_project/src/feature_importance.rs:7:5
26:     0x1002c0550 - wine_quality_project::main::h8ea774ef7312908f
    at /Users/USER/DS210_final_project/wine_quality_project/src/main.rs:29:51
27:     0x1002ad144 - core::ops::function::FnOnce::call_once::h4f9075827d7eb03c
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/ops/function.rs:250:5
28:     0x1002c2ba0 - std::sys::backtrace::_rust_begin_short_backtrace::h589586fe8b7b433d
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/sys/backtrace.rs:152:18
29:     0x1002b3f6c - std::rt::lang_start::{{closure}}::hb75c03ac6c6a457
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/rt.rs:162:18
```

```

30: 0x1002e0bd8 - core::ops::function::impls::<impl core::ops::function::FnOnce<A> for &F>::call_once::h4f74490c6170ea16
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/core/src/ops/function.rs:284:13
31: 0x1002e0bd8 - std::panicking::try::do_call::h2f36d2f1f1af8d28
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:557:40
32: 0x1002e0bd8 - std::panicking::try::ha6af9029a7d93c94
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:521:19
33: 0x1002e0bd8 - std::panic::catch_unwind::ha4f738ae2ba7c3a4
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panic.rs:350:14
34: 0x1002e0bd8 - std::rt::lang_start_internal::{{closure}}::hf216622dc2c733e3
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/rt.rs:141:48
35: 0x1002e0bd8 - std::panicking::try::do_call::haa691957db1dd55f
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:557:40
36: 0x1002e0bd8 - std::panicking::try::ha0c1a49b9fabc98e
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panicking.rs:521:19
37: 0x1002e0bd8 - std::panic::catch_unwind::h68ad032c646bb0ab
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/panic.rs:350:14
38: 0x1002e0bd8 - std::rt::lang_start_internal::hdd117cb81a316264
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/rt.rs:141:20
39: 0x1002b3f38 - std::rt::lang_start::hceb62945b723d8d8
    at /rustc/eeb90cda1969383f56a2637cbd3037bdf598841c/library/std/src/rt.rs:161:17
40: 0x1002c0994 - _main

```

Testing:

Several unit tests ensure the correctness of data ingestion and normalization.

cargo test:

```

Compiling wine_quality_project v0.1.0 (/Users/USER/DS210_final_project/wine_q
uality_project)
Finished `test` profile [unoptimized + debuginfo] target(s) in 0.44s
Running unittests src/main.rs (target/debug/deps/wine_quality_project-a7cf4
2b138a13702)

running 0 tests

test result: ok. 0 passed; 0 failed; 0 ignored; 0 measured; 0 filtered out; fini
shed in 0.00s

```

Interpretation:

The feature importance values are shown as 0.0 for all features in the dataset. This result suggests that the logistic regression weights associated with all features were computed as 0.

This is likely due to the impact of normalization. During feature normalization, if the input values for features have minimal variance or are scaled to a narrow range, the logistic regression algorithm might assign negligible weights to the features.

While the output currently shows zero feature importances and predictions, this provides a baseline for refining the model. The statistical results offer insights into the average normalized alcohol content and quality ratings for red and white wines.

Conclusion

This project laid the groundwork for analyzing and predicting wine quality using Rust. The current implementation establishes a modular and testable framework, with further iterations needed to refine feature importance and prediction accuracy. Rust's efficiency and the dataset's rich attributes ensure a strong potential for meaningful insights into wine quality prediction.

Works Cited

1. Dua, D., & Graff, C. (2019). UCI Machine Learning Repository: Wine Quality Dataset. <https://archive.ics.uci.edu/dataset/186/wine+quality>
2. Rust Programming Language. The Rust Standard Library. <https://doc.rust-lang.org/std/>
3. Towards Data Science. Notes on Graph Theory and Centrality Measurements. <https://towardsdatascience.com/notes-on-graph-theory-centrality-measurements-e37d2e49550a>
4. Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830. (For reference on feature importance and logistic regression concepts).
5. Wikipedia. Logistic Regression. https://en.wikipedia.org/wiki/Logistic_regression