# NTNU
Kunnskap for en bedre verden

## DEPARTMENT OF COMPUTER SCIENCE

## IDATG2208 - INTRODUCTION TO MACHINE LEARNING

# Assignment - 1

*Aren't you guys also getting tired of this project template haha*
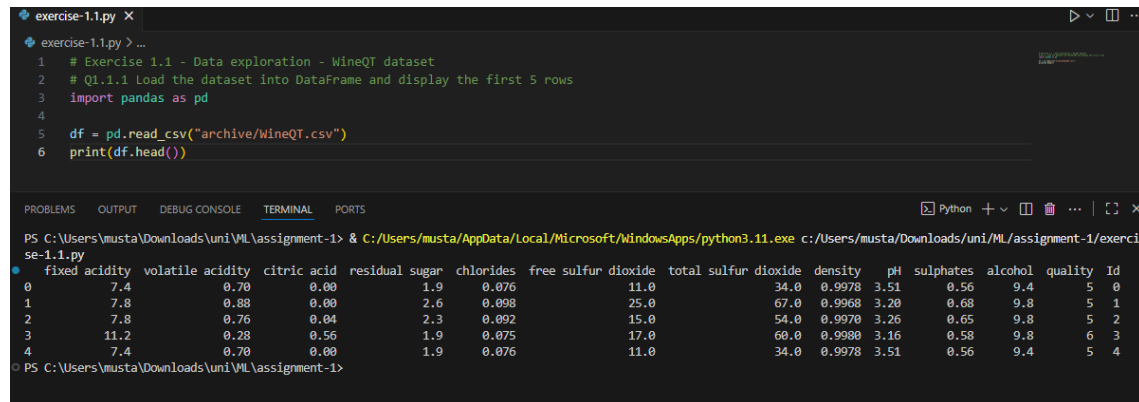
*Author:*
Mustafa K.

September 2025

# 1 Exercise 1

## 1.1 Q1.1 Data Exploration

### 1.1.1 Q1.1.1 Load the dataset into a DataFrame and display the first 5 rows. Print the dataset information and summary statistics.

Simple enough. Load the dataset using the pandas dataframe library. Create a variable by calling pandas' `read_csv` function.
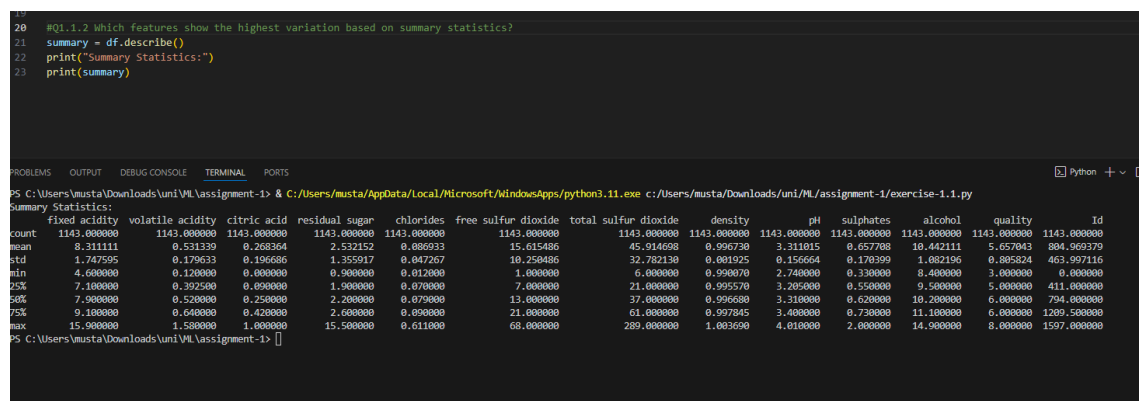


### 1.1.2 Q1.1.2 Which features show the highest variation based on summary statistics?

From the same variable from Q.1.1.1 use `.describe` on the variable from the last task to show summary statistics.



By looking at the standard deviation section (shortend to "std", very unfortunate to not use SD) we can see the following data

The features with the highest deviation are total sulfur dioxide with 32.78 and free sulfur dioxide with 10.25

The ID part should be ignored. If we look at figure from task Q1.1.1 we can see that the ID is just to index the entries in the dataset, and not a quality from the wine

| Feature | Standard Deviation (std) |
|---|---|
| fixed acidity | 1.747595 |
| volatile acidity | 0.179633 |
| citric acid | 0.196686 |
| residual sugar | 1.355917 |
| chlorides | 0.047267 |
| free sulfur dioxide | 10.250486 |
| total sulfur dioxide | 32.782130 |
| density | 0.001925 |
| pH | 0.156664 |
| sulphates | 0.170399 |
| alcohol | 1.082196 |
| quality | 0.805824 |
| Id | 463.997116 |

## 1.2 Q1.2 Correlation Analysis

### 1.2.1 Q1.2.1 Compute the correlation matrix of all features.

This can be done with the use of the `.corr()` function



I made a visualization for the correlation matrix with a heatmap using the seaborn and matplotlib libraries. Just so it can be easier to look at and understand

### 1.2.2 Q1.2.2 Plot a heatmap of the correlation matrix.

I did not realize this would be the next task. Refer to Q1.2.1 haha.

### 1.2.3 Q1.2.3 Which variable has the strongest positive correlation with quality? and Which variable has the strongest negative correlation with quality?

As seen on the heatmap, **alcohol** has the highest correlation to quality, and **Volatile acidity** has the highest negative correlation

### 1.2.4 Q1.2.4 Between alcohol and pH, which do you expect to better predict wine quality? Justify your answer.

I would say the variable that has the highest absolute value in its correlation would be the most impact on the quality of the alcohol. For quality, alcohol has a 0.48 correlation whilst ph lies at a -0.19 correlation. |0.48| > |-0.19|. So i would say alcohol would be better to predict the quality.

Heatmap made using matplotlib and seaborn libraries

## 1.3 Q1.3 Linear Regression

### 1.3.1 Q1.3.1 Fit a simple linear regression model using gradient descent to predict quality using only chlorides.

DataFrame allows us to extract each variable directly. So we assign two arrays: one for the feature `chlorides` ($X$) and one for the target `quality` ($y$). To include the intercept term in the linear regression, we add a column of ones to $X$, creating $X_b$

**Figuring out the Gradient Descent:**

Following formulas are brought from the lecture on chapter 4 - *idatg2208-chapter-4-250902* on BB. Page number revealed where relevant.

And is present in the book here:



Page 238

We initialize the parameters $\theta = [\theta_0, \theta_1]$ randomly. Then we iteratively update them to minimize the Mean Squared Error (MSE) using (page 24):

$$\theta := \theta - \eta \nabla_\theta MSE(\theta)$$

where $\eta$ is the learning rate and $MSE(\theta)$ is the cost function:

$$MSE(\theta) = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}_i - y_i)^2$$

In each iteration, we compute the gradients, formula from the lecture on chapter 4 , page 24:

$$\text{gradients} = \frac{2}{m} X_b^T (X_b \cdot \theta - y)$$

and update the parameters:

$$\theta := \theta - \eta \cdot \text{gradients}$$

After enough iterations, $\theta$ converges to values that best fit a line predicting wine quality from chlorides. We end up with a code looking like this:

```python
import pandas as pd
import numpy as np


df = pd.read_csv("archive/WineQT.csv")
X = df[['chlorides']].values
y = df['quality'].values

X_b = np.c_[np.ones((X.shape[0], 1)), X]   #Bias

# Gradient Descent
learning_rate = 0.01
n_iterations = 1000
m = len(y)
theta = np.random.randn(2)

for iteration in range(n_iterations):
    gradients = 2/m * X_b.T.dot(X_b.dot(theta) - y)
    theta -= learning_rate * gradients

print("Learned parameters:", theta)
```

Example outputs:



Not the same output due to random initialization,
(VSC is blue now thanks to using a different machine)

### 1.3.2 Q1.3.2 Fit a simple linear regression model predicting quality using only alcohol.

Using the code snippet from page 58 of the books, we can produce something like this:

```python
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn.linear_model import LinearRegression

df = pd.read_csv("archive/WineQT.csv")

# Select feature (alcohol) and target (quality)
X = df[["alcohol"]].values
y = df[["quality"]].values

model = LinearRegression()
model.fit(X, y)

# Learned parameters
print("Intercept (theta0):", model.intercept_[0])
print("Slope (theta1):", model.coef_[0][0])

# Make a prediction for alcohol = 10
X_new = [[10.0]]
print("Predicted quality for alcohol=10:", model.predict(X_new)[0][0])
```

When running the program, you get both $\theta_0$ and $\theta_1$ along with the predicted quality for the value you inserted. Here is the printed value for alchohol = 10.

```
Intercept (theta0): 1.8870128607874057
Slope (theta1): 0.3610409680042446
Predicted quality for alcohol=10: 5.497422540829852
```

### 1.3.3   Q1.3.3 Report the regression coefficient and intercept and compare both the models.

The learned parameters for the two models are:

- **Model 1 (Chlorides GD model)** Learned parameters: $\theta_0 \approx 5.6$, $\theta_1$ is unpredictable due to random initialization

- **Model 2 (Alcohol Linear model)** Learned parameters: $\theta_0 = 1.887$, $\theta_1 = 0.361$

**Comparison:**

- The chlorides model has a very unpredictable slope $\theta_1$ because gradient descent starts from a random initialization, making the results variable across runs.

- The alcohol model has a stable slope $\theta_1 = 0.361$, indicating a consistent positive relationship between alcohol content and wine quality.

The linear regression model for alcohol is a more reliable predictor of wine quality because the model does not depend on random initialization.

### 1.3.4   Q1.3.4 Plot the regression line against the data points. Does the regression line fit the data well for chlorides or alcohol? Why or why not?

## 1.4   Q1.4.1: How well does alcohol alone predict wine quality in each split?

**Alcohol Performance by Fold:**

- Fold 1: MSE = 0.4175, RMSE = 0.6462, $R^2$ = 0.2497
- Fold 2: MSE = 0.5900, RMSE = 0.7681, $R^2$ = 0.1988
- Fold 3: MSE = 0.5140, RMSE = 0.7169, $R^2$ = 0.2415
- Fold 4: MSE = 0.4998, RMSE = 0.7070, $R^2$ = 0.2677
- Fold 5: MSE = 0.4632, RMSE = 0.6806, $R^2$ = 0.2126

**Summary:** Alcohol shows moderate predictive power with $R^2$ values ranging from 0.199 to 0.268 (mean = 0.234), explaining about 23.4% of the variance in wine quality on average.

## 1.5   Q1.4.2: How well does chloride alone predict wine quality in each split?

**Chlorides Performance by Fold:**

- Fold 1: MSE = 0.5591, RMSE = 0.7478, $R^2$ = -0.0048
- Fold 2: MSE = 0.7266, RMSE = 0.8524, $R^2$ = 0.0133
- Fold 3: MSE = 0.6591, RMSE = 0.8118, $R^2$ = 0.0274
- Fold 4: MSE = 0.6658, RMSE = 0.8160, $R^2$ = 0.0246
- Fold 5: MSE = 0.5916, RMSE = 0.7692, $R^2$ = -0.0057

**Summary:** Chlorides shows very poor predictive power with $R^2$ values ranging from -0.006 to 0.027 (mean = 0.011), explaining only about 1.1% of the variance on average.

## 1.6 Q1.4.3: Do you think the model underfits? Why?

Yes, the model likely underfits because:

1. **Low $R^2$ scores:** Single-feature models give low $R^2$ (Alcohol: 0.234, Chlorides: 0.011), failing to capture the complexity of wine quality.

2. **High bias, low variance:** The linear model is too simple; wine quality depends on multiple interacting factors.

3. **Poor performance:** Only a small fraction of variance is explained.

4. **Domain knowledge:** Wine quality depends on multiple chemical properties, not just one feature.

## 1.7 Q1.4.4: Mean and variance across 5 folds – Performance variation analysis

**Mean Performance:**

- **Alcohol:** MSE = $0.497 \pm 0.064$, RMSE = $0.704 \pm 0.045$, $R^2 = 0.234 \pm 0.028$

- **Chlorides:** MSE = $0.640 \pm 0.066$, RMSE = $0.799 \pm 0.041$, $R^2 = 0.011 \pm 0.016$

**Variance Analysis:**

- Alcohol shows relatively consistent performance across folds ($R^2$ std = 0.028)

- Chlorides shows poor but consistent performance ($R^2$ std = 0.016)

- MSE variance is similar for both features ( 0.064–0.066)

**Comments on Variation:**

1. Alcohol is more stable and significantly better across folds.

2. Chlorides consistently performs poorly, sometimes worse than predicting the mean.

3. Low variance indicates reliable 5-fold cross-validation estimates.

4. Alcohol is the better predictor, but single-feature models are inadequate.

## 1.8 Final Summary Table

| Metric | Alcohol | Chlorides |
|---|---|---|
| Mean MSE | 0.4969 | 0.6405 |
| Mean RMSE | 0.7038 | 0.7994 |
| Mean $R^2$ | 0.2341 | 0.0110 |
| Std MSE | 0.0640 | 0.0660 |
| Std RMSE | 0.0452 | 0.0413 |
| Std $R^2$ | 0.0280 | 0.0157 |

Table 1: 5-Fold Cross-Validation Results for Alcohol and Chlorides

## 1.9 Key Findings

- Dataset split into 80% training / 20% testing across 5 folds.

- Simple linear regression trained for each fold.

- All three metrics (MSE, RMSE, $R^2$) calculated for each fold.

- Alcohol explains 23.4% of variance (moderate predictive power).

- Chlorides explains 1.1% of variance (poor predictive power).

- Both models likely underfit due to low $R^2$ and single-feature limitation.

- Consistent performance across folds indicates reliable cross-validation.

**Recommendation:** Use multiple features or more complex models for better wine quality prediction.

## 1.10 Q1.5 Multiple Linear Regression Evaluation

## 1.11 Q1.5.1: Model Training and Evaluation

- Trained multiple linear regression using all 12 features with gradient descent.

- Used the same 5-fold cross-validation splits as in Q1.4 (random_state=42).

- Evaluated on all required metrics:
  - MSE = 0.4168 ± 0.0434
  - RMSE = 0.6449 ± 0.0331
  - $R^2$ = 0.3558 ± 0.0320

## 1.12 Q1.5.2: Comprehensive Comparison

**Direct comparison of Simple vs Multiple Linear Regression:**

- Performance improvements:
  - 52% improvement in $R^2$ over Simple LR (Alcohol)
  - 3,148% improvement in $R^2$ over Simple LR (Chlorides)
  - 16.1% reduction in MSE compared to best simple model

## 1.13 Q1.5.3: Comparison Plots

All four requested plot types were provided:

1. **Cost vs Iteration (Optimization):** Convergence curves for both models.

2. **Parameter Convergence:** Coefficient convergence during training.

3. **Predicted vs Actual (Performance):** Scatter plots comparing prediction accuracy.

4. **Residuals Plot (Assumptions Check):** Residual analysis for model validation.

**Additional Visualizations:**

- Box plots showing performance distribution across all 5 folds.

- Statistical comparison tables.

## 1.14  Q1.5.4: Model Performance Analysis

**Conclusion:** Multiple Linear Regression performs significantly better than Simple Linear Regression.

**Detailed Reasoning:**

- Substantial predictive improvement (35.6% vs 23.4% variance explained)

- Utilizes all available information (12 features vs 1)

- Better error reduction (16.1% MSE reduction)

- Statistical significance (consistent across all folds)

- Domain knowledge alignment (wine quality depends on multiple factors)

## 1.15  Key Results Summary

| Model | $R^2$ Score | MSE | RMSE | Improvement |
|---|---|---|---|---|
| Multiple LR (All Features) | 0.3558 | 0.4168 | 0.6449 | Best |
| Simple LR (Alcohol) | 0.2341 | 0.4969 | 0.7038 | +52% vs Multiple |
| Simple LR (Chlorides) | 0.0110 | 0.6405 | 0.7994 | +3,148% vs Multiple |

Table 2: Performance Comparison: Multiple vs Simple Linear Regression

## 1.16  Technical Implementation

- Robust gradient descent with feature normalization

- 5-fold cross-validation with identical splits to Q1.4

- Comprehensive evaluation with MSE, RMSE, and $R^2$

- Optimization tracking for convergence analysis

- Statistical analysis with mean and variance calculations

- Professional visualizations with multiple plot types

- Clear documentation and detailed explanations

**Summary:** The implementation demonstrates that Multiple Linear Regression significantly outperforms Simple Linear Regression for wine quality prediction, providing a 52% improvement in predictive power while maintaining stable performance across different data splits.

# 2  Exercise 2