# Task 02: Model Training Fundamentals with Wine Quality Dataset

Ammad Waqas
Phone: +92 3185875628
InternID: ARCH-2506-0228

July 7, 2025

## 1   Project Overview

This report addresses Task 02, focusing on mastering machine learning algorithms from Chapter 4 of *Hands-On Machine Learning with Scikit-Learn, TensorFlow, and Keras*. The task involves studying linear models, optimization methods, and model evaluation techniques, followed by implementing a regression pipeline on the Wine Quality dataset from the UCI ML Repository. The objectives are to preprocess data, engineer features, train and compare models (Linear Regression, SGDRegressor, Ridge, Lasso, and Elastic Net), plot learning curves, tune hyperparameters, and analyze results. The report includes Chapter 4 exercise solutions, dataset analysis, visualizations, and comparative notes on gradient descent variants and regularization effects.

## 2   Chapter 4 Notes

Chapter 4 covers training machine learning models with a focus on:

- **Linear Models**: Linear Regression uses Ordinary Least Squares (OLS) via the Normal Equation or SVD. Polynomial Regression extends this with feature expansion, increasing overfitting risk. Regularization (Ridge, Lasso, Elastic Net) constrains weights to prevent overfitting.

- **Optimization Methods**: Batch Gradient Descent (GD) uses the full dataset, Stochastic GD uses one instance, and Mini-batch GD uses small batches. Learning rate schedules and convergence diagnostics (e.g., early stopping) improve optimization.

- **Logistic Regression**: Estimates probabilities for binary/multinomial classification using log loss, with linear decision boundaries.

- **Model Evaluation**: Learning curves diagnose bias/variance; hyperparameter tuning optimizes performance.

### 2.1   Comparative Notes

- **Gradient Descent Variants**:
    - *Batch GD*: Converges to the global minimum for convex functions (e.g., MSE) but is slow for large datasets ($O(m)$).

- *Stochastic GD*: Faster, suitable for large datasets, but bounces around the minimum due to randomness. Learning schedules help convergence.
- *Mini-batch GD*: Balances speed and stability, leverages GPU optimization, less erratic than Stochastic GD.
- **Regularization Effects**: Ridge ($\ell_2$) shrinks weights, Lasso ($\ell_1$) sets some to zero (sparse models), and Elastic Net combines both, robust for correlated features.

# 3 Chapter 4 Exercise Solutions

## 3.1 Exercise 1: Algorithm for Millions of Features

Stochastic or Mini-batch Gradient Descent is suitable due to lower computational complexity ($O(n)$) compared to the Normal Equation ($O(n^{2.4})$ to $O(n^3)$).

## 3.2 Exercise 2: Different Feature Scales

Batch GD, Stochastic GD, and regularized models (Ridge, Lasso) suffer from varying feature scales, slowing convergence. Use `StandardScaler` to normalize features.

## 3.3 Exercise 3: Local Minima in Logistic Regression

No, the log loss is convex, ensuring Gradient Descent converges to the global minimum.

## 3.4 Exercise 4: Convergence of Gradient Descent

Batch GD converges to the global minimum for convex functions. Stochastic and Mini-batch GD may not settle precisely due to randomness but can converge with a proper learning schedule.

## 3.5 Exercise 5: Increasing Validation Error in Batch GD

Increasing validation error indicates overfitting. Fix by reducing model complexity, adding regularization, or using early stopping.

## 3.6 Exercise 6: Stopping Mini-batch GD on Validation Error Increase

Not ideal, as validation error fluctuates. Wait for consistent increases, then roll back to the best model (early stopping).

## 3.7 Exercise 7: Fastest GD and Convergence

Stochastic GD reaches the vicinity fastest due to frequent updates. Batch GD converges precisely. Use learning schedules for Stochastic/Mini-batch GD convergence.

## 3.8 Exercise 8: Large Gap in Learning Curves

A large gap indicates overfitting. Solutions: increase training data, reduce polynomial degree, or apply regularization (Ridge, Lasso).

### 3.9 Exercise 9: High Training/Validation Error in Ridge

High errors suggest high bias (underfitting). Reduce $\alpha$ to decrease regularization strength.

### 3.10 Exercise 10: Why Use Regularization

- *Ridge vs. Linear*: Ridge prevents overfitting by constraining weights.
- *Lasso vs. Ridge*: Lasso performs feature selection by setting weights to zero.
- *Elastic Net vs. Lasso*: Elastic Net handles correlated features better.

### 3.11 Exercise 11: Outdoor/Indoor and Daytime/Nighttime

Use two Logistic Regression classifiers for independent labels (multioutput).

### 3.12 Exercise 12: Batch GD with Early Stopping for Softmax Regression

```python
import numpy as np
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import log_loss
iris = load_iris()
X, y = iris.data[:, (2, 3)], iris.target
X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2,
    random_state=42)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_val_scaled = scaler.transform(X_val)
m, n = X_train_scaled.shape
K = len(np.unique(y))
theta = np.random.randn(n, K)
eta, epochs, tol = 0.01, 1000, 1e-4
min_val_loss, best_theta, best_epoch = float('inf'), None, None
for epoch in range(epochs):
    scores = X_train_scaled.dot(theta)
    exp_scores = np.exp(scores)
    probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
    y_one_hot = np.zeros((m, K))
    y_one_hot[np.arange(m), y_train] = 1
    gradients = (1/m) * X_train_scaled.T.dot(probs - y_one_hot)
    theta -= eta * gradients
    val_scores = X_val_scaled.dot(theta)
    val_probs = np.exp(val_scores) / np.sum(np.exp(val_scores), axis=1,
        keepdims=True)
    val_loss = log_loss(y_val, val_probs)
    if val_loss < min_val_loss:
        min_val_loss, best_theta, best_epoch = val_loss, epoch
    elif epoch - best_epoch > 50:  # Early stopping
        break
```

# 4  Dataset Description

The Wine Quality dataset (UCI ML Repository) contains 4,898 instances of white wine with 11 physicochemical features (e.g., fixed acidity, alcohol) and a quality score (0–10). The task is regression to predict quality.

# 5  Implementation

The pipeline includes:

- **Preprocessing**: Handled missing values (none in dataset), scaled features using `StandardScaler`.

- **Feature Engineering**: Added polynomial features (degree=2) using `PolynomialFeatures`.

- **Models Trained**: Linear Regression (Normal Equation), SGDRegressor, Ridge, Lasso, Elastic Net.

- **Hyperparameter Tuning**: Grid search on $\alpha$ (Ridge, Lasso, Elastic Net) and learning rate (SGD).

- **Evaluation**: RMSE, $R^2$, training time, and feature importance (coefficient magnitudes).

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler, PolynomialFeatures
from sklearn.linear_model import LinearRegression, Ridge, Lasso, ElasticNet
from sklearn.linear_model import SGDRegressor
from sklearn.metrics import mean_squared_error, r2_score
import matplotlib.pyplot as plt
import time

# Load dataset
data = pd.read_csv('winequality-white.csv', sep=';')
X = data.drop('quality', axis=1)
y = data['quality']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)

# Preprocessing
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
poly = PolynomialFeatures(degree=2, include_bias=False)
X_train_poly = poly.fit_transform(X_train_scaled)
X_test_poly = poly.transform(X_test_scaled)

# Train models
models = {
    'Linear': LinearRegression(),
    'SGD': SGDRegressor(penalty=None, random_state=42),
    'Ridge': Ridge(alpha=1.0),
    'Lasso': Lasso(alpha=0.1),
    'ElasticNet': ElasticNet(alpha=0.1, l1_ratio=0.5)
}
```

```python
results = []
for name, model in models.items():
    start_time = time.time()
    model.fit(X_train_poly, y_train)
    train_time = time.time() - start_time
    y_pred = model.predict(X_test_poly)
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    r2 = r2_score(y_test, y_pred)
    results.append({'Model': name, 'RMSE': rmse, 'R2': r2, 'Training Time':
        train_time})

# Hyperparameter tuning (example for Ridge)
param_grid = {'alpha': [0.01, 0.1, 1, 10]}
grid_search = GridSearchCV(Ridge(), param_grid, cv=5, scoring='
    neg_mean_squared_error')
grid_search.fit(X_train_poly, y_train)
print("Best Ridge alpha:", grid_search.best_params_)

# Learning curves
def plot_learning_curves(model, X, y):
    X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2,
        random_state=42)
    train_errors, val_errors = [], []
    for m in range(1, len(X_train), 100):
        model.fit(X_train[:m], y_train[:m])
        y_train_pred = model.predict(X_train[:m])
        y_val_pred = model.predict(X_val)
        train_errors.append(np.sqrt(mean_squared_error(y_train[:m],
            y_train_pred)))
        val_errors.append(np.sqrt(mean_squared_error(y_val, y_val_pred)))
    plt.plot(train_errors, 'r-+', label='Training RMSE')
    plt.plot(val_errors, 'b-', label='Validation RMSE')
    plt.xlabel('Training Set Size')
    plt.ylabel('RMSE')
    plt.legend()
    plt.savefig(f'learning_curve_{model.__class__.__name__}.png')
    plt.close()

for name, model in models.items():
    plot_learning_curves(model, X_train_poly, y_train)

# Coefficient analysis plot
plt.figure(figsize=(10, 6))
for name, model in models.items():
    if hasattr(model, 'coef_'):
        plt.plot(np.abs(model.coef_), label=name)
plt.xlabel('Feature Index')
plt.ylabel('Coefficient Magnitude')
plt.title('Coefficient Magnitudes Across Models')
plt.legend()
plt.savefig('coefficient_plot.png')
plt.close()
```

| Model | RMSE | $R^2$ | Training Time (s) | Non-zero Coefficients |
|---|---|---|---|---|
| Linear | 0.73 | 0.45 | 0.02 | 77 |
| SGD | 0.74 | 0.44 | 0.05 | 77 |
| Ridge | 0.72 | 0.46 | 0.03 | 77 |
| Lasso | 0.75 | 0.43 | 0.04 | 45 |
| ElasticNet | 0.74 | 0.44 | 0.04 | 60 |

Table 1: Model Performance on Wine Quality Dataset

# 6 Comparative Analysis

## 6.1 Model Comparison

## 6.2 Analysis

- **Performance**: Ridge performed best (lowest RMSE, highest $R^2$) due to regularization reducing overfitting on polynomial features. Lasso had fewer non-zero coefficients, indicating feature selection.

- **Polynomial Degree**: Degree=2 improved performance over linear features but risked overfitting without regularization.

- **Applications**: Predicting wine quality for quality control in wineries.

## 6.3 Coefficient Analysis

Figure 1: Coefficient Magnitudes for Regularized Models

Lasso reduced many coefficients to zero, unlike Ridge, which shrank them uniformly.

# 7 Learning Curve Visualizations

Learning curves (saved as $\text{learning}_c urve_*.png) show$ :

```
Linear/SGD: High RMSE plateau (underfitting).

Ridge/Lasso/ElasticNet: Lower training RMSE but slight gap with validation,
indicating mild overfitting mitigated by regularization.
```

# 8 References

- UCI ML Repository: https://archive.ics.uci.edu/ml/datasets/wine+quality

- Hands-On ML GitHub: https://github.com/ageron/handson-ml2

- Scikit-Learn Linear Models: https://scikit-learn.org/stable/modules/linear_model.html

- 3Blue1Brown Gradient Descent: https://www.youtube.com/watch?v=IHZwWFHWa-w

- StatQuest Regularization: https://www.youtube.com/watch?v=Q81RR3yKn30

# 9 GitHub Repository

All code, data, and visualizations are available at: https://github.com/gxammad/Manual_Task2.git. The repository includes winequality-white.csv, $model_{training.ipynb}$, requir