# Data Processing – Final Project

Iker Aldasoro Marculeta

Jon Lejardi Jericó
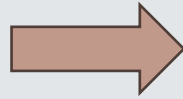
# Table of Contents

# Introduction

# Input Variable Analysis

Average Ratings per Category

Correlation between Numerical Features

Rating Distribution by Year

Distribution of Ratings

Rating Distribution and Outliers

# Text Preprocessing

**Tokenization** → **Homogeneization** → **Cleaning** → **Vectorization**

**Tokenization**
- Divide phrases into words

**Homogeneization**
- Lemmatization
- Convert to lowercase

**Cleaning**
- Remove stop words, extra white spaces, special characters …

**Code:**

```python
def NTLK_clean(text):
    """NLT pipeline using NTLK

    Args:
        text (category): Category of the dataset to apply the NLP

    Returns:
        _type_: Tokenized word
    """
    stop_words = set(stopwords.words("english"))
    if not isinstance(text, str):
        return ""

    try:
        # Initialize lemmatizer
        lemmatizer = WordNetLemmatizer()

        # Convert to lowercase
        text = text.lower()

        # Remove special characters and digits
        text = re.sub(r"[^a-zA-Z\s]", "", text)

        # Remove extra whitespace
        text = " ".join(text.split())

        # Remove extra whitespace
        text = lemmatizer.lemmatize(text)

        # Basic word tokenization (split by space)
        tokens = word_tokenize(text)

        # Remove stopwords
        tokens = [token for token in tokens if token not in stop_words]

        return " ".join(tokens)
    except Exception as e:
        print(f"Error in text preprocessing: {str(e)}")
        return ""
```

**Output:**

```
3. Applying NLT Pipeline
0    1. Place the stock, lentils, celery, carrot, t...
1    Combine first 9 ingredients in heavy medium sa...
2    In a large heavy saucepan cook diced fennel an...
3    Heat oil in heavy large skillet over medium-hi...
4    Preheat oven to 350°F. Lightly grease 8x8x2-in...
Name: directions_pre, dtype: object

0    place stock lentils celery carrot thyme salt m...
1    combine first ingredients heavy medium saucepa...
2    large heavy saucepan cook diced fennel onion b...
3    heat oil heavy large skillet mediumhigh heat a...
4    preheat oven f lightly grease xxinch glass bak...
Name: directions_post, dtype: object
```

# Vector Representation

Convert text data into **Word vectors** (numerical) to make inputs compatible to ML models.

- <u>TF-IDF</u>: Context-agnostic and focuses on word importance (20130, 1000).

- <u>Word2Vec</u>: Captures shallow word relationships, lacks full contextual understanding: (20130, 100), vector size 100 and window 5. -> Mean vector

- <u>BERT</u>: Context-aware and provides deep contextual embeddings (20130, 768).
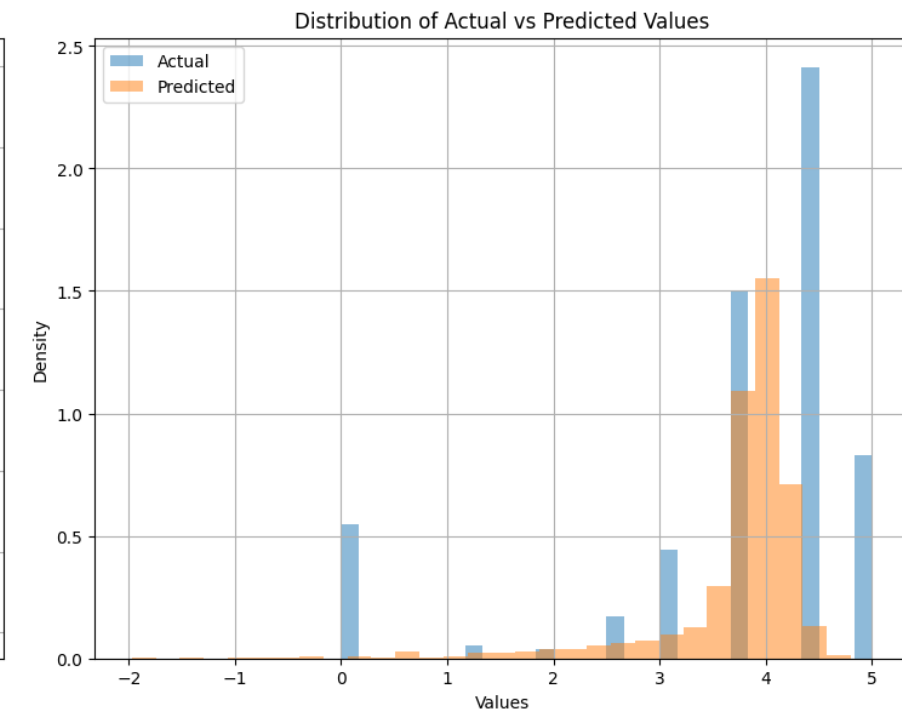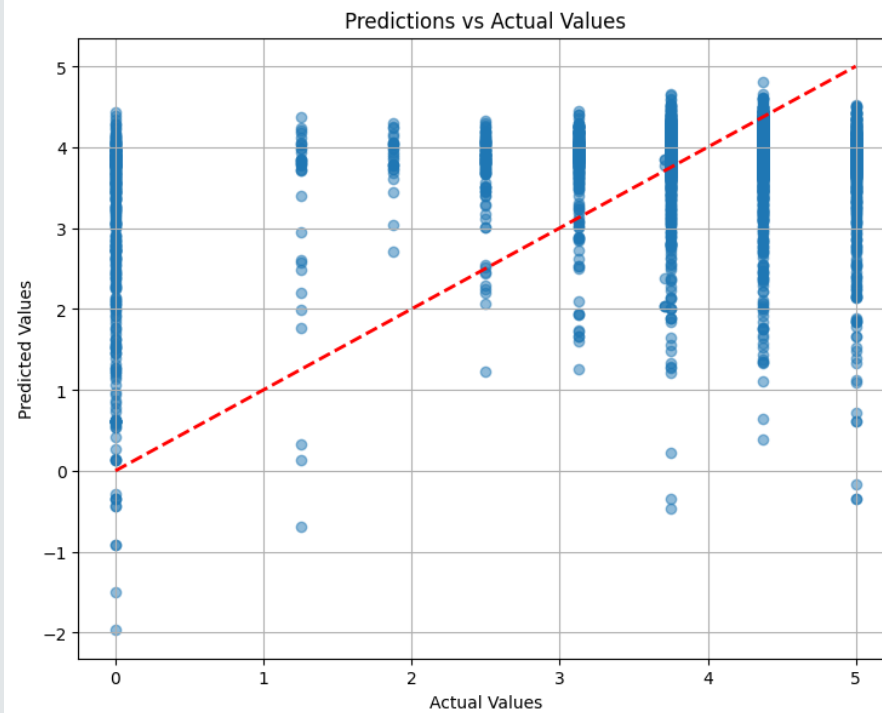
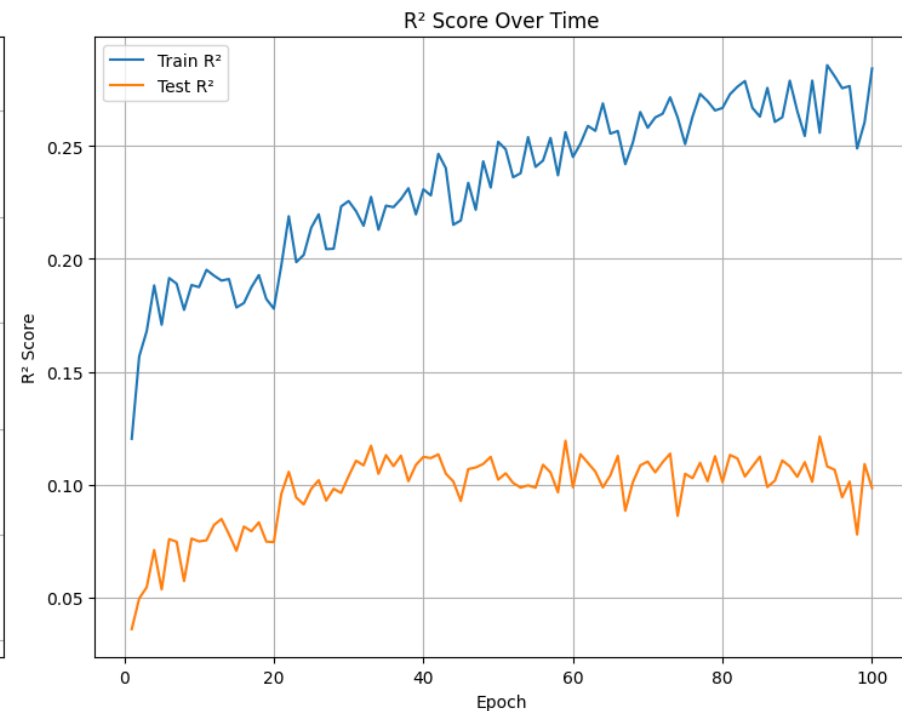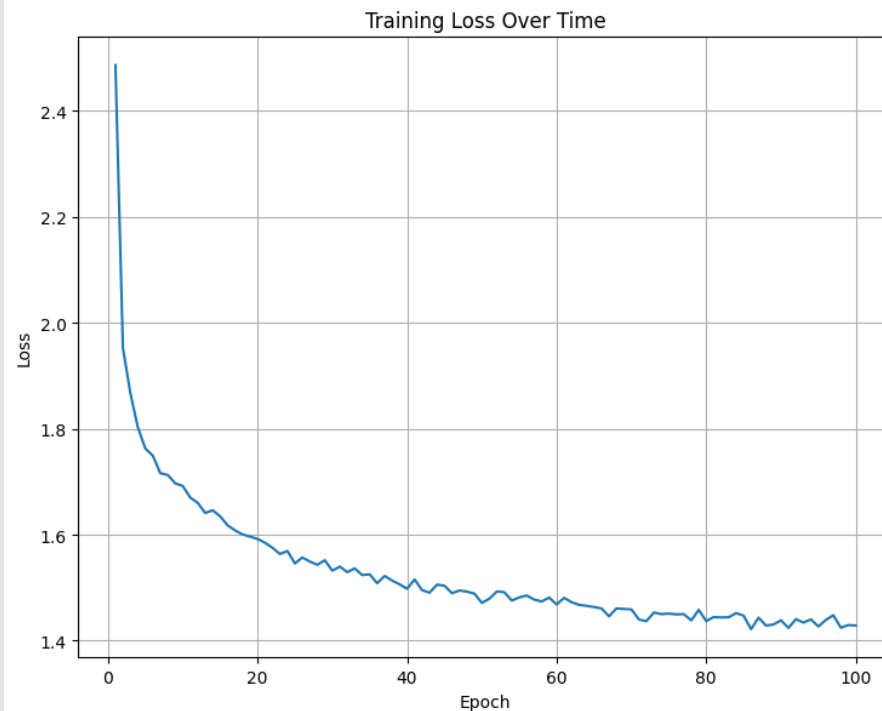# Regression Models

# Neural Networks

Multi-Layer NN for regression

- 2 hidden layers (128 and 64 features).
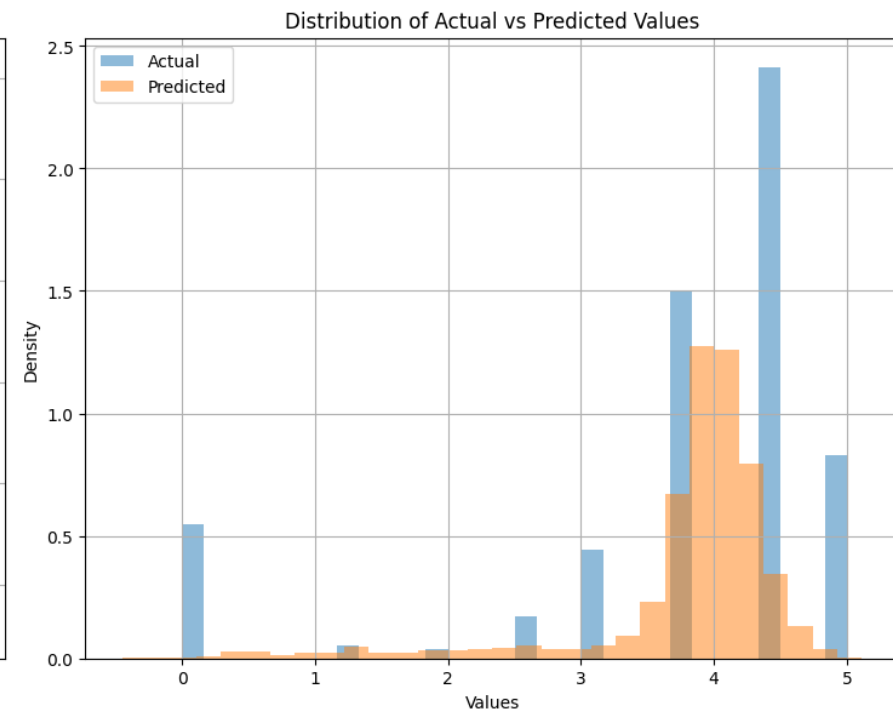
- ReLu activation function.
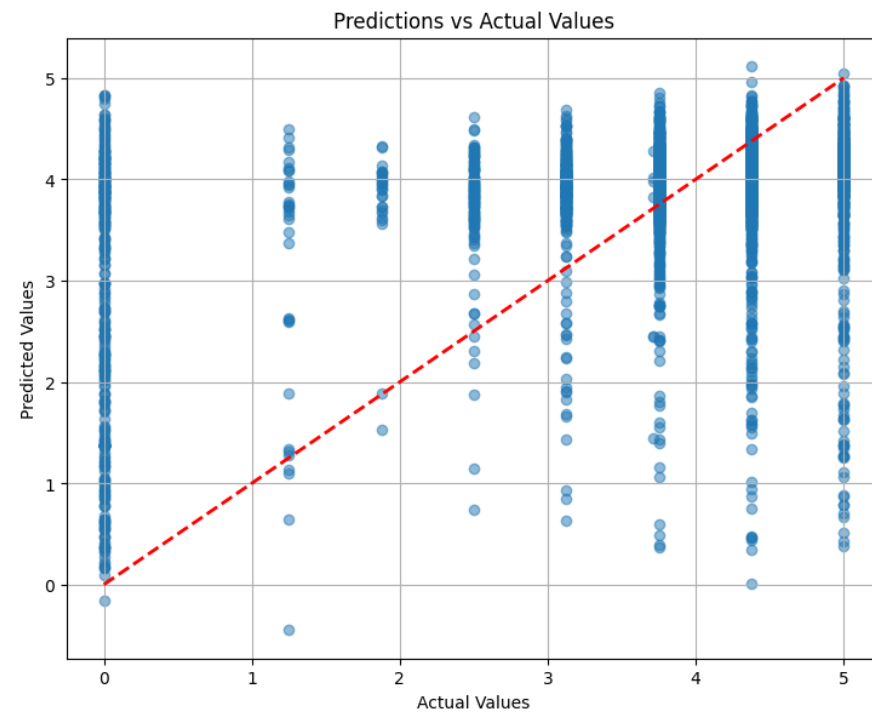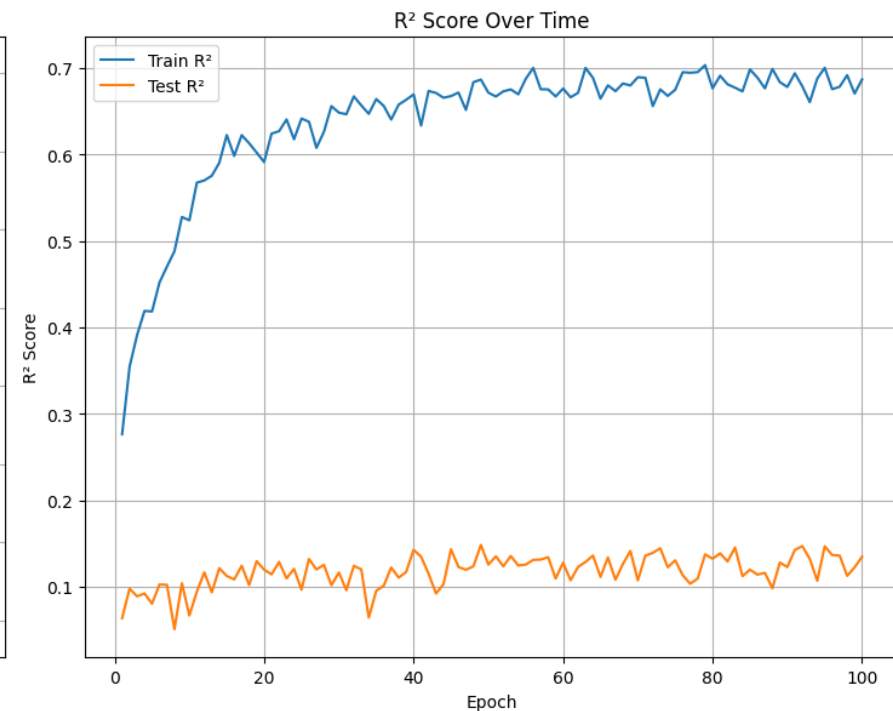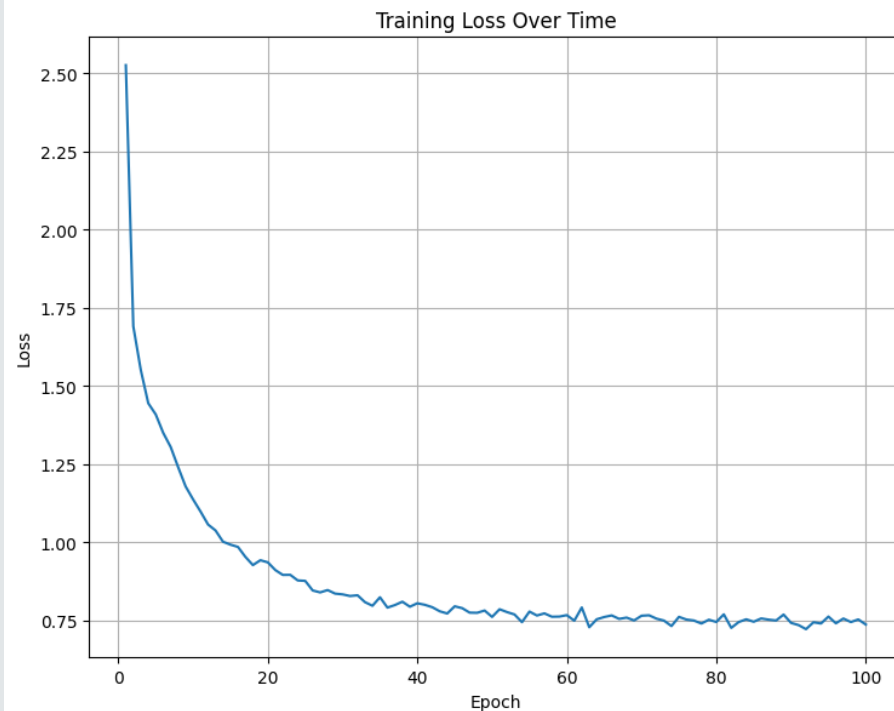
- 30% dropout.

```python
model_ML = nn.Sequential(
    nn.Linear(input_size, 128),
    nn.ReLU(),
    nn.Dropout(0.3),   # Add dropout
    nn.Linear(128, 64),
    nn.ReLU(),
    nn.Dropout(0.3),   # Add dropout
    nn.Linear(64, 1),
)
```

|  | Train $R^2$ | Test $R^2$ |
|---|---|---|
| **TD-IDF** | 0,6869 | 0,1346 |
| **Word2Vec** | 0,206 | 0,1339 |
| **BERT** | 0,2843 | 0,0985 |

**TF-IDF**

Training Loss Over Time

R² Score Over Time

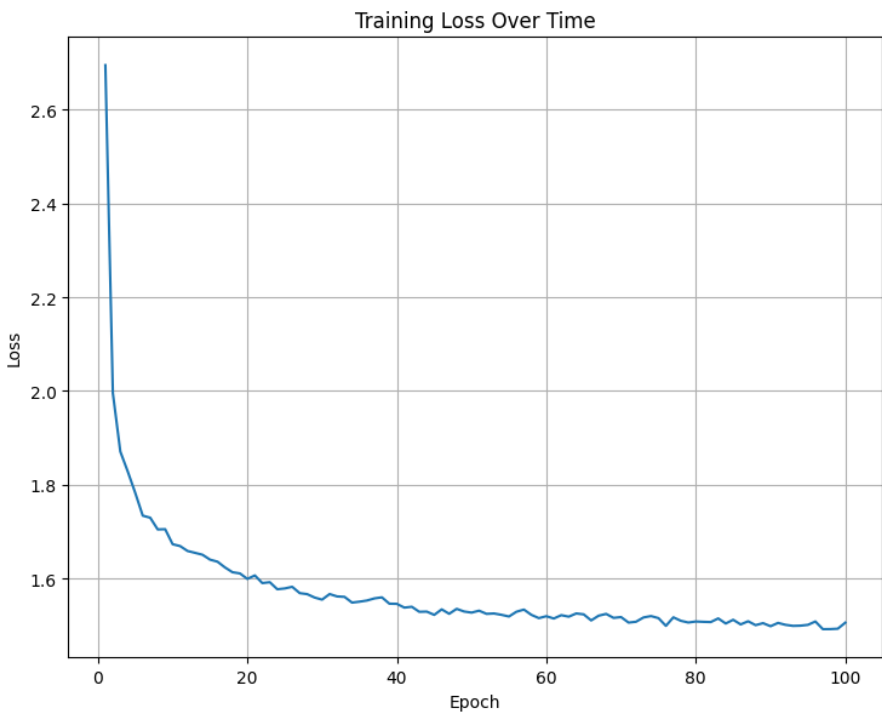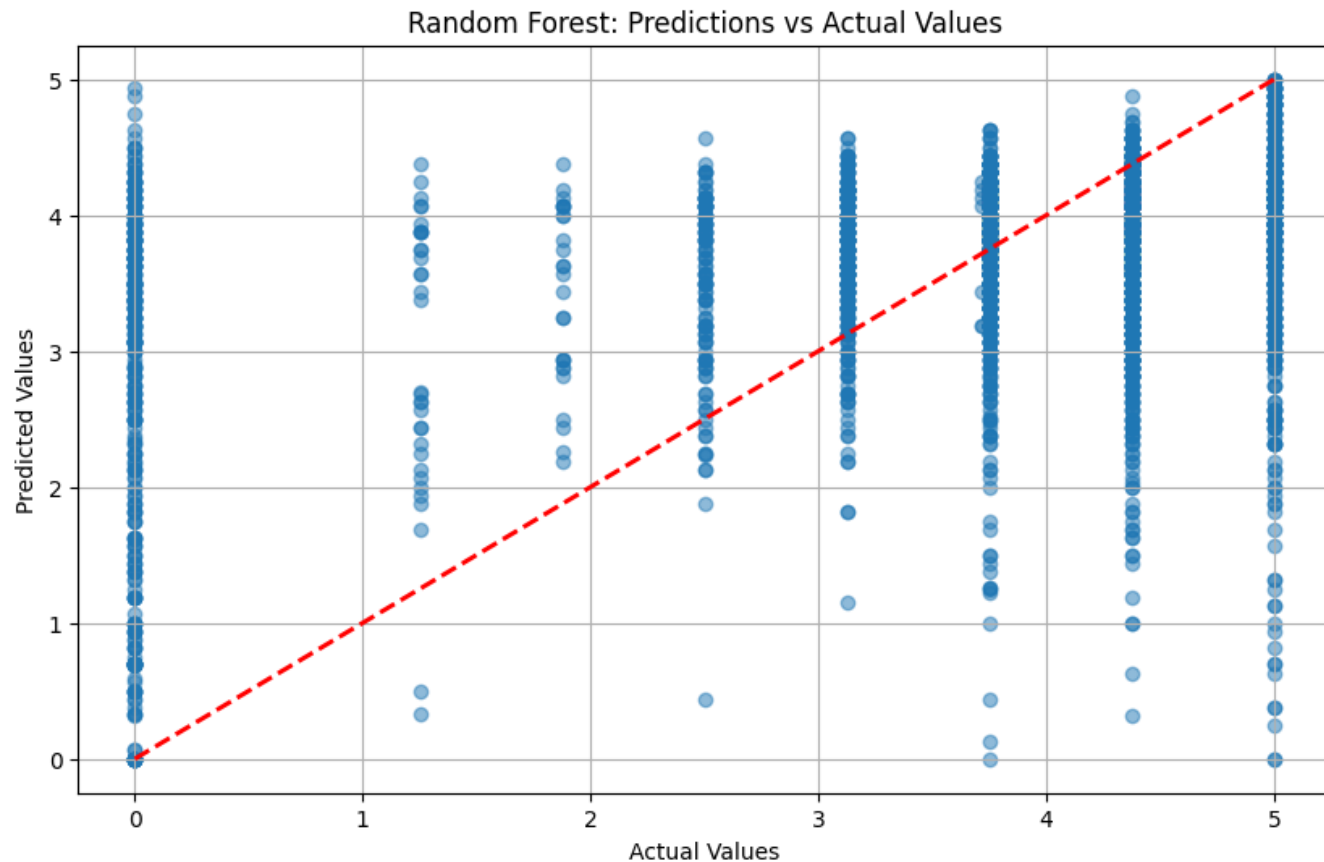Predictions vs Actual Values

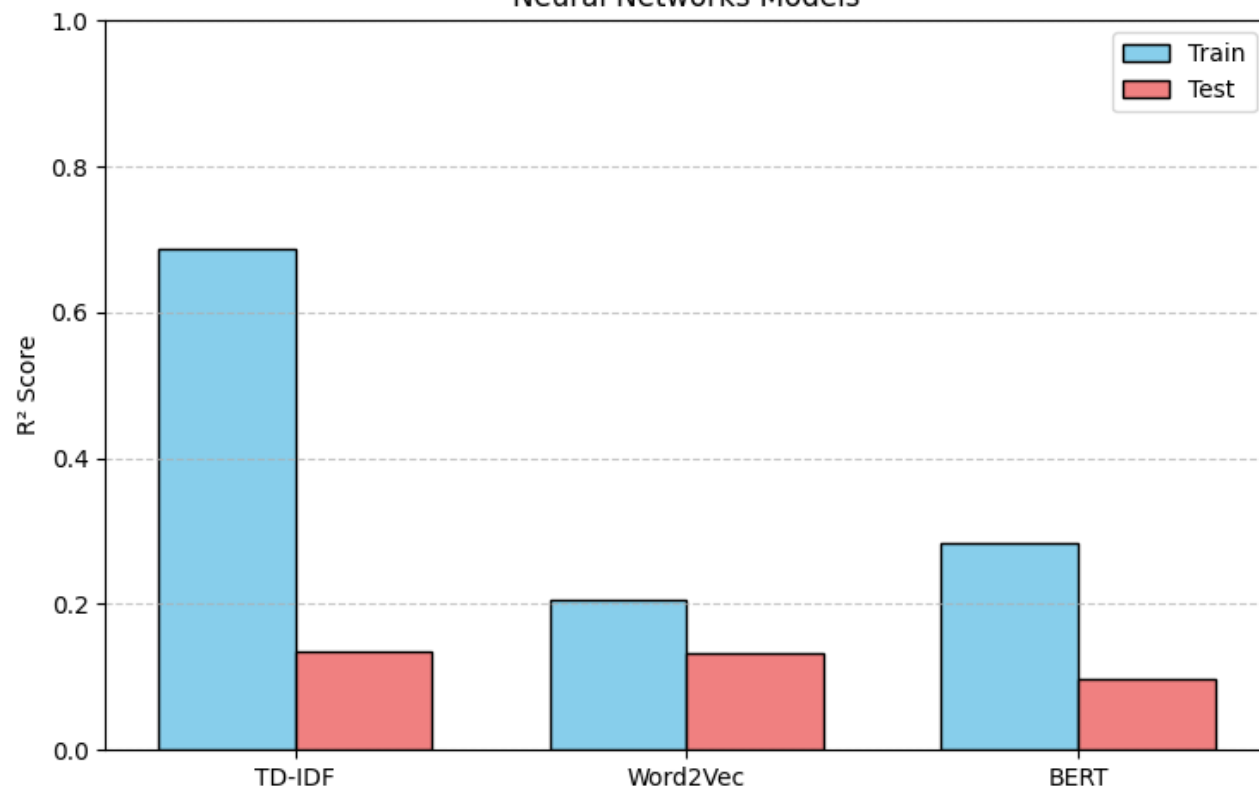Distribution of Actual vs Predicted Values

# Random Forest



```python
# Train Random Forest
rf_model = RandomForestRegressor(n_estimators=10, random_state=42)
# rf_model.fit(X_train_t.numpy(), y_train_t.numpy().ravel())
rf_model.fit(X_train, y_train)

# Get predictions
rf_train_predictions = rf_model.predict(X_train)
rf_test_predictions = rf_model.predict(X_test)

# Calculate metrics
rf_train_r2 = r2_score(y_train, rf_train_predictions)
rf_test_r2 = r2_score(y_test, rf_test_predictions)
rf_train_mse = mean_squared_error(y_train, rf_train_predictions)
rf_test_mse = mean_squared_error(y_test, rf_test_predictions)
```
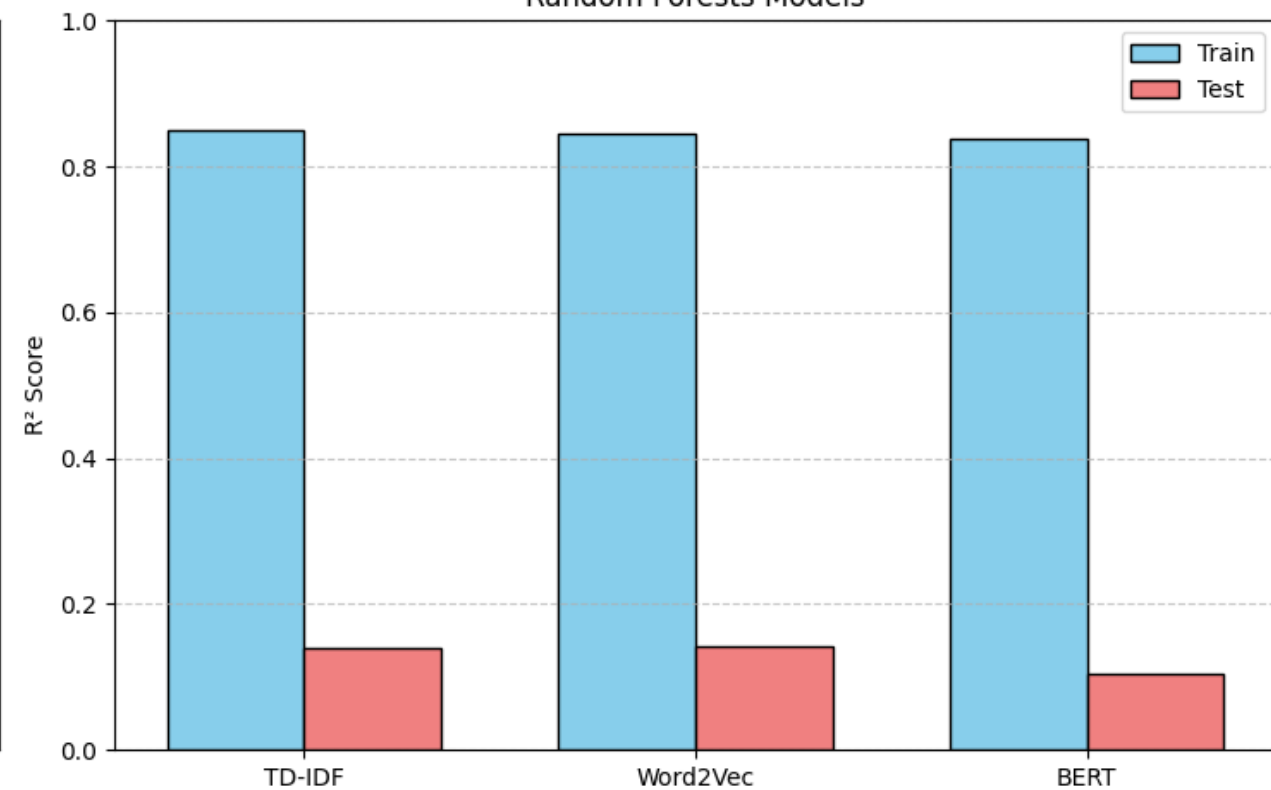
Random Forest: Predictions vs Actual Values

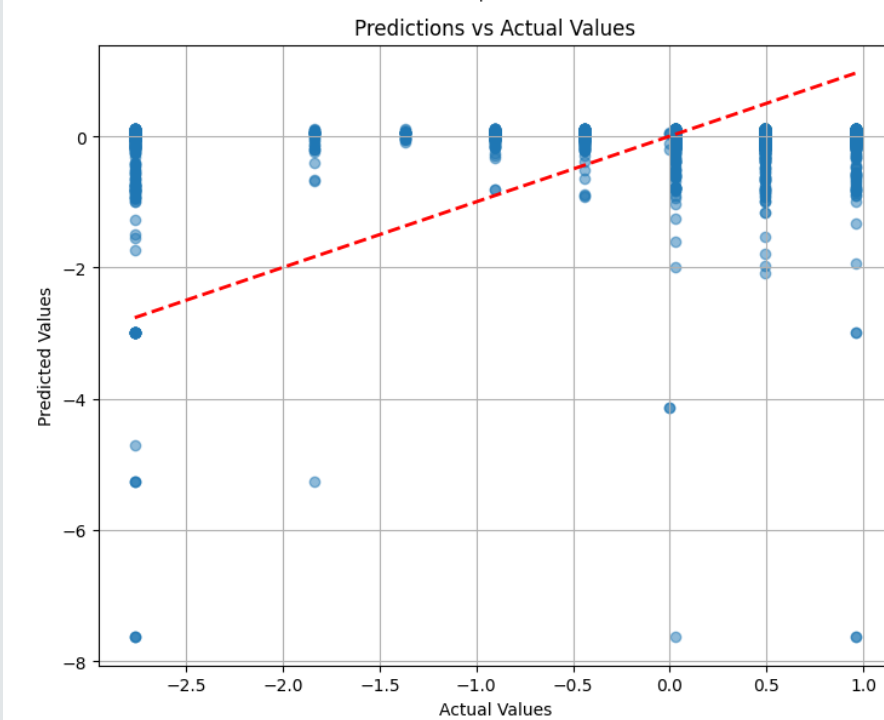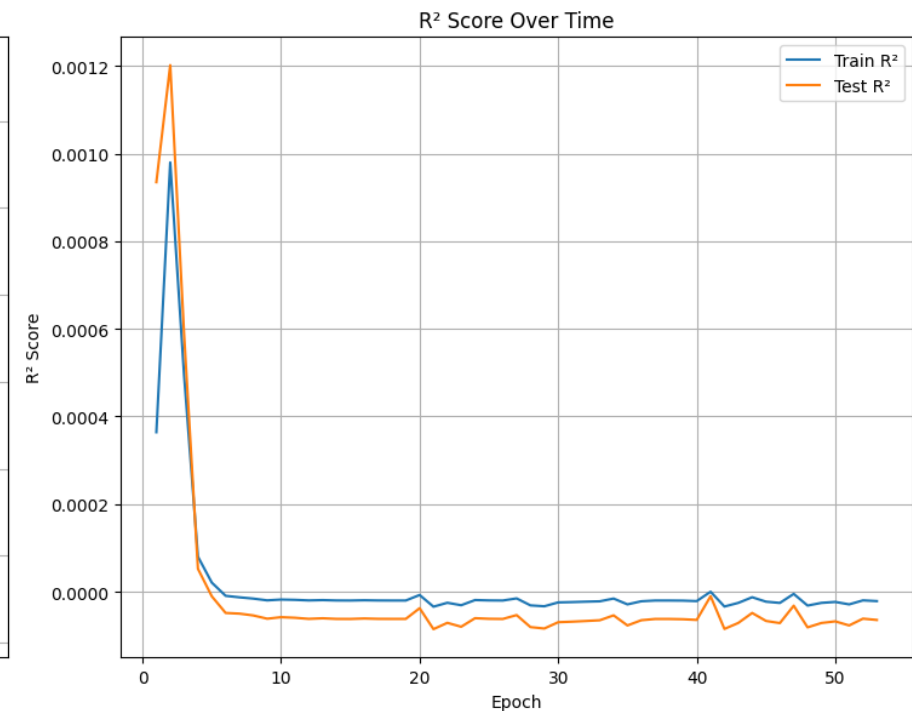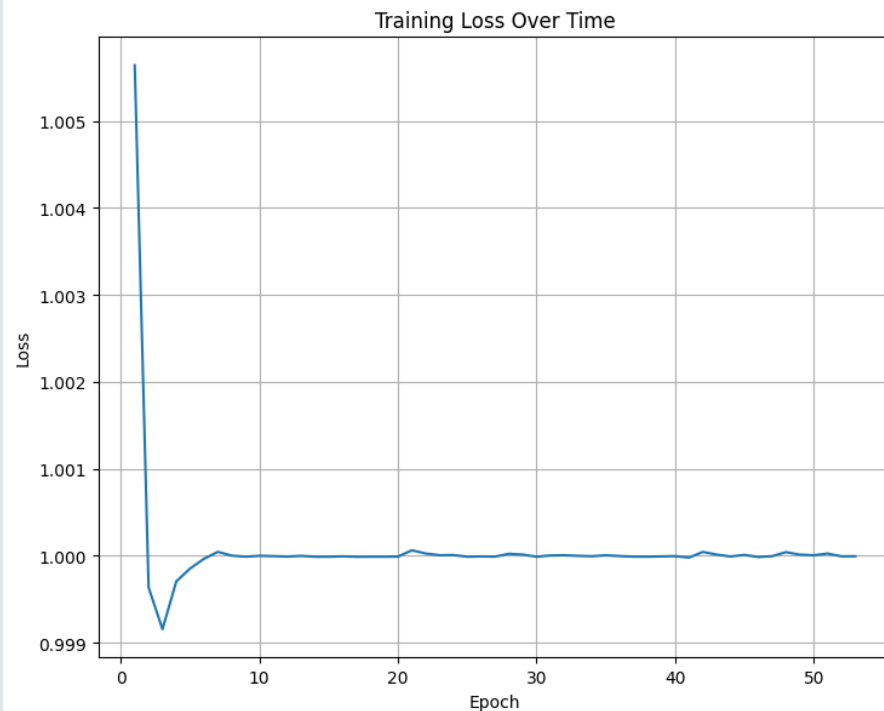|  | Train $R^2$ | Test $R^2$ |
|---|---|---|
| **TD-IDF** | 0,8489 | 0,1388 |
| **Word2Vec** | 0,844 | 0,1419 |
| **BERT** | 0,8378 | 0,1037 |

# Hugging Face

We used RoBERTa as a model but the results were not better

# Conclusion

## Problems

- Bad regression model results.

## Possible solutions

- Try different models.
- Find correlated inputs.

# THANK YOU VERY MUCH

Iker Aldasoro

Jon Lejardi