# Report: Assignment #4

ITE4005, **Data Science.**
2016025305, **Jihun Kim**<jihunkim@hanyang.ac.kr**>.**

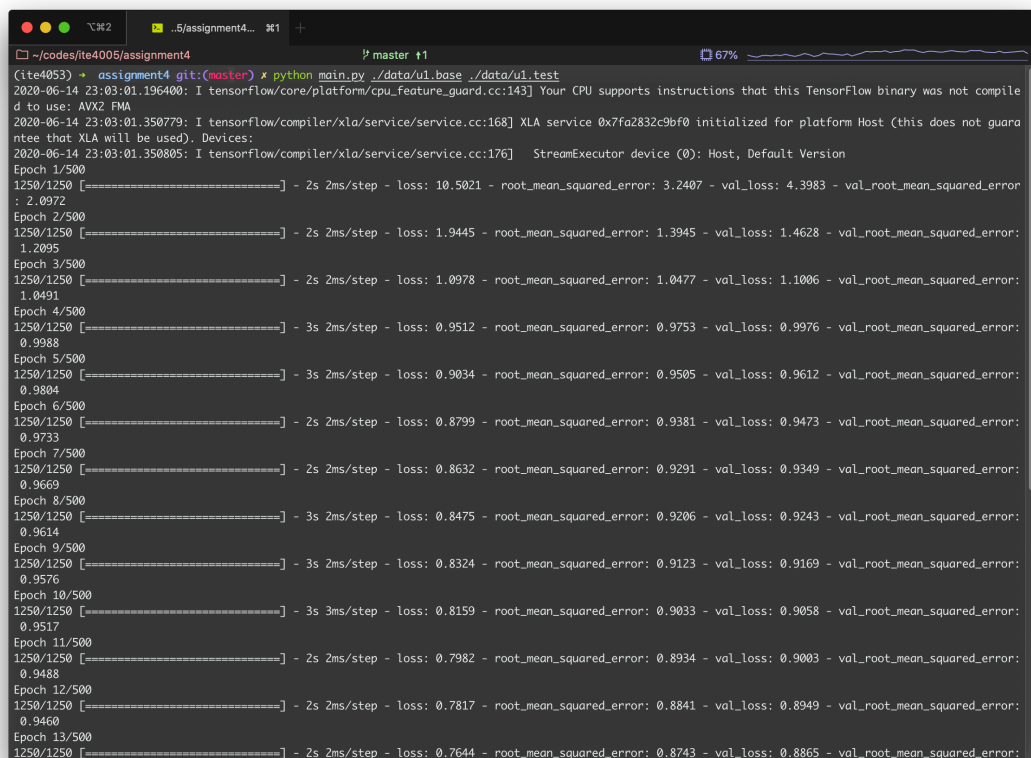Predict the ratings of movies in test data by using the given training data containing movie ratings of users.

# Getting Started

## Development Environment

*   **OS:** macOS 10.15.5
*   **Language:** Python 3.8.3 (TensorFlow 2.2.0, NumPy 1.18.5)

## Run

```
$ cd /path/to/repo/assignment4
$ pip install -r requirements.txt
$ python main.py ./data/u1.base ./data/u1.test
```

# Implementation

For predicting the ratings, I used simple correlative filtering. Implementation has done with Keras with TensorFlow as backend.

## Model Construction

I constructed simple model that performs dot product of two embedded vectors. One vector represents a user, and another vector represents a movie. The dimension of vectors are empirically set for 20. Lower dimensions occurred the increase of loss, while higher dimensions occurred overfitting.

```python
def cf_model(n_users, n_movies):
    user_input = layers.Input(shape=(1,))
    user_x = layers.Embedding(n_users, 20, input_length=1, name='user_embed')(user_input)
    item_input = layers.Input(shape=(1,))
    item_x = layers.Embedding(n_movies, 20, input_length=1, name='item_embed')(item_input)
    rating = layers.Dot(axes=-1)([user_x, item_x])
    rating = layers.Flatten()(rating)
    return models.Model([user_input, item_input], rating)
```

## Model Compilation

... and I compiled model. Used Adam as optimizer, because it shows great performance in most situations.

```python
model.compile(optimizer=optimizers.Adam(1e-3),
              loss=losses.mean_squared_error,
              metrics=[metrics.RootMeanSquaredError()])
```

## Training Model

Next step is training, of course. User and movie vector embeddings are trained in training phase. Adopted early stopping method in training phase.

```python
model.fit([train_data['user_id'], train_data['item_id']], train_data['rating'],
          validation_data=([test_data['user_id'], test_data['item_id']], test_data['rating']),
          callbacks=callbacks.EarlyStopping(patience=2),
          batch_size=64, epochs=500, verbose=1)
```
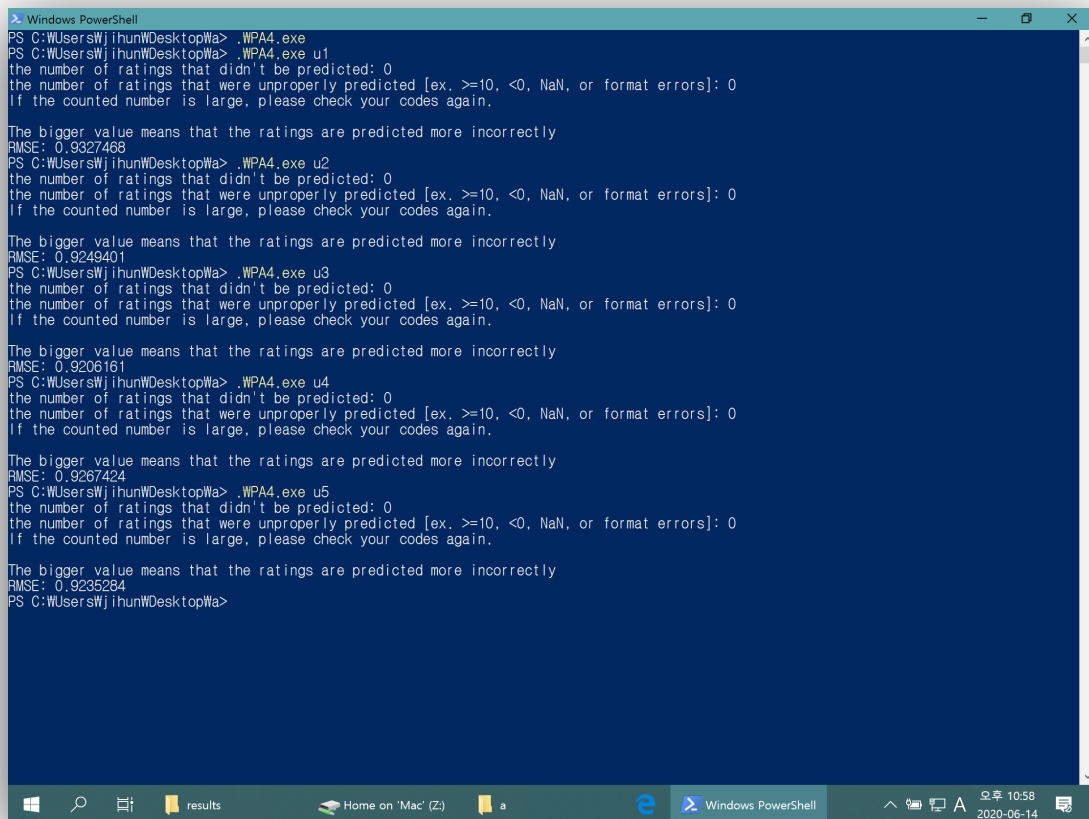
## Evaluate on Test Set

For evaluation, I constructed new model which does exactly same thing with previouly constructed model, except it has single additional layer which make network output fit in range [1, 5].

```python
rating = layers.Lambda(lambda x : K.minimum(K.maximum(x, 1), 5))(model.output)
final_model = models.Model(model.input, rating)
final_model.compile(loss=losses.mean_squared_error, metrics=[metrics.RootMeanSquaredError()])
final_model.evaluate([test_data['user_id'], test_data['item_id']], test_data['rating'])
    test_data['rating'] = final_model.predict([test_data['user_id'], test_data['item_id']])
```

# Results

Results are below. For all train/test sets, my model shows 0.92 ~ 0.94 RMSE with no errors.