# Recommender Systems

In this lab, we'll be using Keras to build a recommender system. We'll be using the MovieLens dataset, a common benchmark dataset for recommender systems.

MovieLens is a web-based recommender system and virtual community that recommends movies for its users to watch, based on their film preferences using collaborative filtering of members' movie ratings and movie reviews. You can check out the website here: https://movielens.org/

We will download a subset of the dataset containing 100k ratings. There are tens of millions of ratings in the full dataset, spanning hundreds of thousands of users and movies. The subset we'll be using is a good example to demonstrate the concepts in this lab.

In [ ]:
```python
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
from zipfile import ZipFile
from urllib.request import urlretrieve
from pathlib import Path

ML_100K_URL = "http://files.grouplens.org/datasets/movielens/ml-100k.zip"
ML_100K_FILENAME = Path("ml-100k.zip")
ML_100K_FOLDER = Path("ml-100k")

if not ML_100K_FOLDER.exists():
    if not ML_100K_FILENAME.exists():
        urlretrieve(ML_100K_URL, ML_100K_FILENAME.name)
    with ZipFile(ML_100K_FILENAME.name) as zip_file:
        zip_file.extractall()
```

The goal of MovieLens is to enable models to predict the rating a user would give to a movie they have not yet watched. This is a classic example of a recommendation system. The dataset is huge, and contains many parts giving information about the movies, the users, and the ratings. To begin with, we will look at the ratings file. Each line in the ratings file (u.data) is formatted as:

```
user_id, item_id, rating, timestamp
```

Which tells us a single user's rating of a single movie.

We will start by loading the ratings data into a pandas dataframe and then take a look at the first few rows. If you haven't used Pandas before, it's an extremely

powerful library for dealing with tabular data. You can think of it as a Python version of Excel.

```python
import pandas as pd

raw_ratings = pd.read_csv(ML_100K_FOLDER / "u.data", sep='\t',
                          names=["user_id", "item_id", "rating", "timestamp"
raw_ratings
```

The second file we'll look at is the movie metadata. This file (u.item) contains information about each movie, including the title and release date. Each line in the file is formatted as:

```
movie_id | movie_title | release_date | video_release_date |
IMDb_URL | unknown | Action | Adventure | Animation | Children's |
Comedy | Crime | Documentary | Drama | Fantasy | Film-Noir | Horror
| Musical | Mystery | Romance | Sci-Fi | Thriller | War | Western
```

As you can see, the genres are binary variables. As with one-hot encoding, a 1 indicates that the movie is of that genre, and a 0 indicates that it is not. We aren't going to work with the genre data in this lab, but it's easy to imagine that it could be useful in a real-world recommendation system.

```python
columns_to_keep = ['item_id', 'title', 'release_date', 'video_release_date',
items = pd.read_csv(ML_100K_FOLDER / "u.item", sep='|', names=columns_to_kee
                    encoding='latin-1', usecols=range(5))
items
```

By default, the `release_date` column is a string. We can convert it to a `datetime` object using the `pd.to_datetime` function. This will make it easier to work with in the future (if we want to do things like check which date came first, for example).

We can also extract the year from the date and store it in a separate column. This will make it easier to do things like plot the number of movies released each year.

```python
items['release_date'] = pd.to_datetime(items['release_date']) # Pandas makes
items['release_year'] = items['release_date'].dt.year # For later use
```

For our purposes, it will be easier to work with the data if we merge our two dataframes into a single dataframe. We can do this using the `merge` method. We'll merge the `items` dataframe into the `raw_ratings` dataframe, using the `item_id` column as the key. This will add the movie title and release year to each rating.

```python
all_ratings = pd.merge(items, raw_ratings)
```

```
In [ ]: all_ratings.head()
```

## Data preprocessing

It's always important to understand the data you've collected. Thankfully, Pandas continues to make this easy for us. Using the `describe` method, we can get a quick statistical summary of the data.

```
In [ ]: all_ratings.describe()
```

Let's do a bit more pandas magic to compute the popularity of each movie (number of ratings). We will use the `groupby` method to group the dataframe by the `item_id` column and then use the `size` method to compute the number of ratings for each movie. We will use the `reset_index` method to convert the resulting Series into a dataframe with an `item_id` column.

```
In [ ]: popularity = all_ratings.groupby('item_id').size().reset_index(name='popular
        items = pd.merge(popularity, items)
```

```
In [ ]: items['popularity'].plot.hist(bins=30);
```

```
In [ ]: (items['popularity'] == 1).sum() # Number of movies with only one rating
```

```
In [ ]: items.nlargest(10, 'popularity')['title'] # Get the 10 most popular movies
```

```
In [ ]: all_ratings = pd.merge(popularity, all_ratings)
        all_ratings.describe()
```

```
In [ ]: indexed_items = items.set_index('item_id')
```

```
In [ ]: all_ratings.head()
```

**Quick Exercise**:

As we have seen, the `groupby` method is a powerful tool to quickly compute statistics on the data. Use it to compute the average rating for each movie.

**Hint**: you can use the `mean` method after the `groupby` method.

```
In [ ]: raise NotImplementedError("Please calculate the average rating for each movi
```

Let's split the enriched data in a train / test split to make it possible to do predictive modeling:

```
In [ ]: from sklearn.model_selection import train_test_split

        ratings_train, ratings_test = train_test_split(
```

```
    all_ratings, test_size=0.2, random_state=0)

user_id_train = np.array(ratings_train['user_id'])
item_id_train = np.array(ratings_train['item_id'])
rating_train = np.array(ratings_train['rating'])

user_id_test = np.array(ratings_test['user_id'])
item_id_test = np.array(ratings_test['item_id'])
rating_test = np.array(ratings_test['rating'])
```
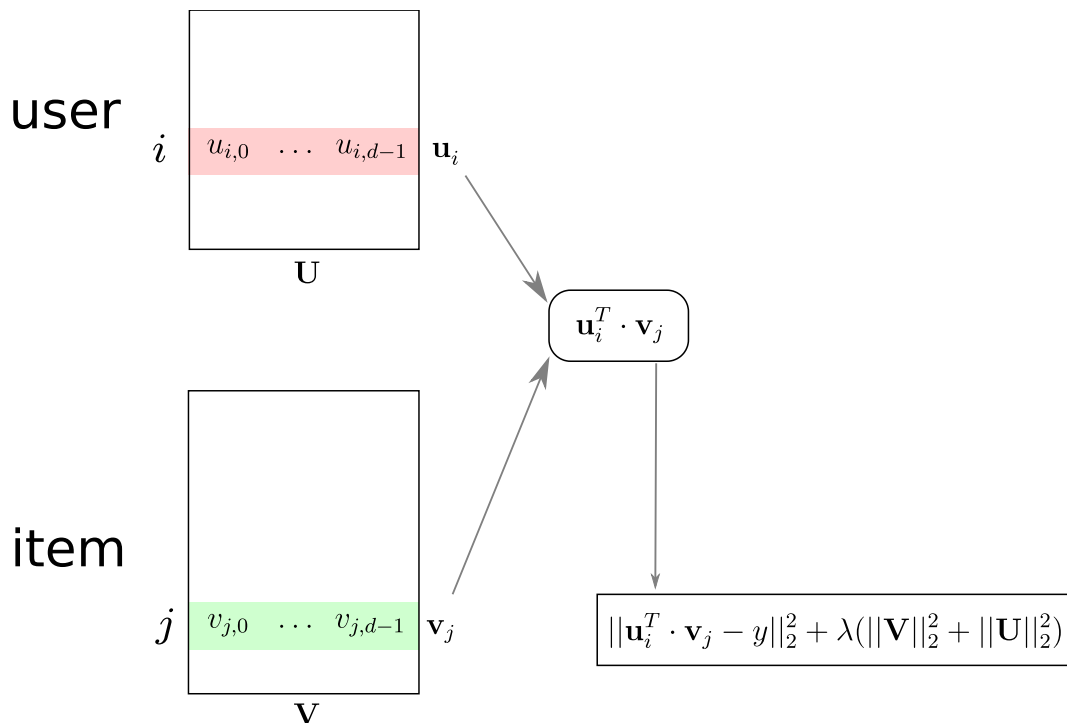
# Explicit feedback: supervised ratings prediction

Now let's begin to do some recommendation! We will build a model that takes a user and a movie as input and outputs a predicted rating. We will be taking advantage of embeddings to represent users and movies. That means that each movie and user will have an abstract representation in a continuous vector space. The model will learn these representations based on the ratings.

## Predictive ratings as a regression problem

The following code implements the following architecture:



```
In [ ]:  from tensorflow.keras.layers import Embedding, Flatten, Dense, Dropout
         from tensorflow.keras.layers import Dot
         from tensorflow.keras.models import Model
```

```
In [ ]:  # For each sample we input the integer identifiers
         # of a single user and a single item
         class RegressionModel(Model):
             def __init__(self, embedding_size, max_user_id, max_item_id):
                 super().__init__()

                 self.user_embedding = Embedding(output_dim=embedding_size,
                                                 input_dim=max_user_id + 1,
                                                 name='user_embedding')
                 self.item_embedding = Embedding(output_dim=embedding_size,
                                                 input_dim=max_item_id + 1,
                                                 name='item_embedding')

                 # The following two layers don't have parameters.
                 self.flatten = Flatten()
                 self.dot = Dot(axes=1)

             def call(self, inputs):
                 user_inputs = inputs[0]
                 item_inputs = inputs[1]

                 user_vecs = self.flatten(self.user_embedding(user_inputs))
                 item_vecs = self.flatten(self.item_embedding(item_inputs))

                 y = self.dot([user_vecs, item_vecs])
                 return y


         model = RegressionModel(embedding_size=64, max_user_id=all_ratings['user_id'
         model.compile(optimizer="adam", loss='mae')
```

## Monitoring runs

When training a model with Keras, we get a `history` object back that contains lots of information about the training run. We can use this to plot the training and validation loss to see how the model has improved during training.

```
In [ ]:  %%time

         # Training the model
         history = model.fit([user_id_train, item_id_train], rating_train,
                             batch_size=64, epochs=10, validation_split=0.1,
                             shuffle=True)
```

```
In [ ]:  plt.plot(history.history['loss'], label='train')
         plt.plot(history.history['val_loss'], label='validation')
         plt.ylim(0, 2)
         plt.legend(loc='best')
         plt.title('Loss');
```

**Questions**:

- Does it look like our model has overfit? Why or why not?

  Your Answer: _____
- Suggest something we could do to prevent overfitting.

  Your Answer: _____

Now that the model is trained, let's check out the quality of predictions:

```python
In [ ]: def plot_predictions(y_true, y_pred):
            plt.figure(figsize=(4, 4))
            plt.xlim(-1, 6)
            plt.xlabel("True rating")
            plt.ylim(-1, 6)
            plt.ylabel("Predicted rating")
            plt.scatter(y_true, y_pred, s=60, alpha=0.01)
```

```python
In [ ]: from sklearn.metrics import mean_squared_error
        from sklearn.metrics import mean_absolute_error

        test_preds = model.predict([user_id_test, item_id_test])
        print("Final test MSE: %0.3f" % mean_squared_error(test_preds, rating_test))
        print("Final test MAE: %0.3f" % mean_absolute_error(test_preds, rating_test)
        plot_predictions(rating_test, test_preds)
```

This graph shows us the range of predicted ratings our model gives, organized by the true rating. We can see that generally, the higher the true rating the higher the predicted rating, although there is quite a range of predictions for each instance. That's okay - our model is very simple, and human preferences are very complex!

Taking a look at the Mean Absolute Error, hopefully you got something around 0.75. This means that, on average, our predicted ratings are about 0.75 stars off from the true ratings. This is a pretty good result for a first attempt. We could probably do better with a more complex model, but we'll leave that for another time.

## Model Embeddings

Our model was built with two embedding layers. These layers have learned a representation of both the users and the movies in our dataset. We can extract these representations and use them to find similar movies or users. We can also do interesting exploratory analysis, like finding the most popular movies among our users, or finding the users that are most interested in a given movie.

```python
In [ ]: # weights and shape
        weights = model.get_weights()
        [w.shape for w in weights]
```

```python
user_embeddings = weights[0]
item_embeddings = weights[1]
```

```python
item_id = 181
print(f"Title for item_id={item_id}: {indexed_items['title'][item_id]}")
```

```python
print(f"Embedding vector for item_id={item_id}")
print(item_embeddings[item_id])
print("shape:", item_embeddings[item_id].shape)
```

As we discussed in lecture, our embeddings are not directly interpretable - we can't look at, say, a value of 0.297 in the embedding vector and say "this means that the movie is a drama". As an aside, there is a field of research dedicated to making *interpretable* embeddings, but it's not something we'll cover in this course.

## Finding our most similar items

Now we can have some fun, investigating the embeddings we've learned. We can start by finding the most similar items to a given item. We can do this by computing the cosine similarity between the item's embedding and the embedding of every other item. We can use the `cosine_similarity` function from `sklearn` to do this.

```python
from sklearn.metrics.pairwise import cosine_similarity

def cosine(a, b):
    # Reshape to the shape our function expects
    a = a.reshape(1, -1)
    b = b.reshape(1, -1)
    return cosine_similarity(a, b)[0, 0]
```

```python
def print_similarity(item_a, item_b, item_embeddings, titles):
    print(titles[item_a])
    print(titles[item_b])
    similarity = cosine(item_embeddings[item_a],
                        item_embeddings[item_b])
    print(f"Cosine similarity: {similarity:.3}")

print_similarity(50, 181, item_embeddings, indexed_items["title"])
```

It makes sense that the original Star Wars, and its later sequel Return of the Jedi have a high similarity. Let's try some other examples:

```python
print_similarity(181, 288, item_embeddings, indexed_items["title"])
```

```python
print_similarity(181, 1, item_embeddings, indexed_items["title"])
```

```
In [ ]:  print_similarity(181, 181, item_embeddings, indexed_items["title"])
```

*Quick Exercise*:

- Find some other films and compare their similarity. Do the results make sense to you? Can you find a pair of films that are very *dissimilar*?

```
In [ ]:  # Code to help you search for a movie title
         partial_title = "Jedi"
         indexed_items[indexed_items['title'].str.contains(partial_title)]

         raise NotImplementedError("Please implement the next steps yourself")
```

Sometimes, even without knowing anything about a user, we can recommend films by asking them about a film that they do like. The code below compares the similarity of a given film to all others, and returns the most similar films.

```
In [ ]:  def most_similar(item_id, item_embeddings, titles,
                          top_n=30):
             # Compute the cosine similarity between the item and all other items
             sims = cosine_similarity(item_embeddings[item_id].reshape(1, -1),
                                      item_embeddings).ravel()

             # [::-1] makes it possible to reverse the order of a numpy
             # array, this is required because most similar items have
             # a larger cosine similarity value
             sorted_indexes = np.argsort(sims)[::-1]
             idxs = sorted_indexes[0:top_n]
             return list(zip(idxs, titles[idxs], sims[idxs]))

         # Find the most similar films to "Star Wars"
         most_similar(50, item_embeddings, indexed_items["title"], top_n=10)
```

```
In [ ]:  # Find the most similar films to "Star Trek VI: The Undiscovered Country"
         most_similar(227, item_embeddings, indexed_items["title"], top_n=10)
```

The similarities do not always make sense: the number of ratings is low and the embedding does not automatically capture semantic relationships in that context. Better representations arise with higher number of ratings, and less overfitting in models or maybe better loss function, such as those based on implicit feedback.

## Visualizing embeddings using TSNE

The t-SNE algorithm enables us to visualize high dimensional vectors in a 2D space by preserving local neighborhoods. We can use it to get a 2D visualization of the item embeddings and see if similar items are close in the embedding space.

```python
In [ ]:  from sklearn.manifold import TSNE

         item_tsne = TSNE(learning_rate="auto", init="pca", perplexity=30).fit_transf
```

```python
In [ ]:  import matplotlib.pyplot as plt

         plt.figure(figsize=(10, 10))
         plt.scatter(item_tsne[:, 0], item_tsne[:, 1]);
         plt.xticks(()); plt.yticks(());
         plt.show()
```

```python
In [ ]:  import plotly.express as px

         tsne_df = pd.DataFrame(item_tsne, columns=["tsne_1", "tsne_2"])
         tsne_df["item_id"] = np.arange(item_tsne.shape[0])
         tsne_df = tsne_df.merge(items.reset_index())

         px.scatter(tsne_df, x="tsne_1", y="tsne_2",
                    color="popularity",
                    hover_data=["item_id", "title", "popularity"])
```

## Exercises

- Add another layer to the neural network and retrain, compare train/test error.
- Try adding more dropout and change layer sizes.

## A recommendation function for a given user

Once the model is trained, the system can be used to recommend a few items for a user that they haven't seen before. The following code does that.

- we use the `model.predict` to compute the ratings a user would have given to all items
- we build a function that sorts these items and excludes those the user has already seen.

```python
In [ ]:  def recommend(user_id, top_n=10):
             item_ids = range(1, items['item_id'].max())
             seen_mask = all_ratings["user_id"] == user_id
             seen_movies = set(all_ratings[seen_mask]["item_id"])
             item_ids = list(filter(lambda x: x not in seen_movies, item_ids))

             user = np.zeros_like(item_ids)
             user[:len(item_ids)] = user_id
             items_ = np.array(item_ids)
             ratings = model.predict([user, items_]).flatten()
             top_items = ratings.argsort()[-top_n:][::-1]
             return [(indexed_items.loc[item_id]["title"], ratings[item_id]) for item
```

```
In [ ]:   for title, pred_rating in recommend(5):
              print("    %0.1f: %s" % (pred_rating, title))
```

## Exercises

- Try modifying our neural network to improve recommendation. You could try adding more layers, or using a different loss function.
- Your goal is to improve the Mean Absolute Error on the test set. Show the results of your best model.

```
In [ ]:   # Extend and improve the model below
          class RegressionModel(Model):
              def __init__(self, embedding_size, max_user_id, max_item_id):
                  super().__init__()

                  self.user_embedding = Embedding(output_dim=embedding_size,
                                                  input_dim=max_user_id + 1,
                                                  name='user_embedding')
                  self.item_embedding = Embedding(output_dim=embedding_size,
                                                  input_dim=max_item_id + 1,
                                                  name='item_embedding')

                  # The following two layers don't have parameters.
                  self.flatten = Flatten()
                  self.dot = Dot(axes=1)

              def call(self, inputs):
                  user_inputs = inputs[0]
                  item_inputs = inputs[1]

                  user_vecs = self.flatten(self.user_embedding(user_inputs))
                  item_vecs = self.flatten(self.item_embedding(item_inputs))

                  y = self.dot([user_vecs, item_vecs])
                  return y


          model = RegressionModel(embedding_size=64, max_user_id=all_ratings['user_id'
          model.compile(optimizer="adam", loss='mae')
```

```
In [ ]:   # Training the model
          history = model.fit([user_id_train, item_id_train], rating_train,
                              batch_size=64, epochs=10, validation_split=0.1,
                              shuffle=True)
```

```
In [ ]:
```