

Lab 1 Exercise - Playing with gradients and matrices in PyTorch

Jonathon Hare (jsh2@ecs.soton.ac.uk)

February 7, 2024

This is the first of a series of exercises that you need to work through **on your own** after completing the accompanying lab session. You'll need to write up your results/answers/findings and submit this to ECS handin as a PDF document along with the other lab exercises near the end of the module (1 pdf document per lab). If you are asked to implement a function, please supply a code snippet in your file.

You should use *no more* than one side of A4 to cover your responses to *this* exercise. This exercise is worth 5% of your overall module grade.

1 Implement a matrix factorisation using gradient descent

You saw (briefly) in the lectures that a form of low-rank matrix factorisation can be achieved by considering the following optimisation problem:

$$\min_{\hat{U}, \hat{V}} (\|A - \hat{U}\hat{V}^\top\|_F^2) \quad (1)$$

where $A \in \mathbb{R}^{m \times n}$, $\hat{U} \in \mathbb{R}^{m \times r}$, $\hat{V} \in \mathbb{R}^{n \times r}$ and $r < \min(m, n)$. The following pseudocode gives an algorithm that uses gradient descent to optimise this problem (note that $\hat{U}_{r,*}$ refers to the entire r -th row of \hat{U}):

Inputs:

A : $\mathbb{R}^{m \times n}$ input matrix.
 r : rank of the factorisation.
 N : number of epochs.
 η : learning rate.

Outputs:

\hat{U} : $\mathbb{R}^{m \times r}$ matrix initialised with values from $\mathcal{U}(0,1)$.
 \hat{V} : $\mathbb{R}^{n \times r}$ matrix initialised with values from $\mathcal{U}(0,1)$.

Algorithm:

```
For epoch = 1 to N
  For r = 1 to m
    For c = 1 to n
       $e = A_{rc} - \hat{U}_{r,*}(\hat{V}_{c,*})^\top$ 
       $\hat{U}_{r,*} \leftarrow \hat{U}_{r,*} + \eta \cdot e \cdot \hat{V}_{c,*}$ 
       $\hat{V}_{c,*} \leftarrow \hat{V}_{c,*} + \eta \cdot e \cdot \hat{U}_{r,*}$ 
return  $\hat{U}, \hat{V}$ 
```

1.1 Implement gradient-based factorisation (1 mark)

Implement the above code method using PyTorch by completing the following method:

```
from typing import Tuple
```

```
def sgd_factorise(A: torch.Tensor, rank: int, num_epochs=1000, lr=0.01) -> Tuple[torch.Tensor, torch.Tensor]:
```

Include the code in your report.

1.2 Factorise and compute reconstruction error (1 mark)

Use your code to compute the rank-2 factorisation of the following matrix and state the result (use the default values for learning rate and number of epochs):

$$\begin{bmatrix} 0.3374 & 0.6005 & 0.1735 \\ 3.3359 & 0.0492 & 1.8374 \\ 2.9407 & 0.5301 & 2.2620 \end{bmatrix}$$

Compute the value of the reconstruction loss^a $\|\mathbf{A} - \hat{\mathbf{U}}\hat{\mathbf{V}}^\top\|_{\mathbb{F}}^2$.

^aHint: you can use `torch.nn.functional.mse_loss` with `reduction='sum'` to do this.

2 Compare your result to truncated SVD

A truncated Singular Value Decomposition is defined as:

$$\tilde{\mathbf{A}} = \mathbf{U}_t \mathbf{\Sigma}_t \mathbf{V}_t^\top \quad (2)$$

where $\tilde{\mathbf{A}} \in \mathbb{R}^{m \times n}$, $\mathbf{U} \in \mathbb{R}^{m \times t}$, $\mathbf{S} \in \mathbb{R}^{t \times t}$, $\mathbf{V} \in \mathbb{R}^{n \times t}$ and $t < \min(m, n)$.

2.1 Compare to the truncated-SVD (1 mark)

Compute the SVD of the matrix \mathbf{A} using `torch.svd`. Set the last singular value to zero, and compute the reconstruction. What is the error? How does it compare to the result from your algorithm above? Why is this the case^a?

^aMaybe google Eckart-Young theorem

3 Matrix completion

One of the really neat things we can do with our gradient-based approach to matrix factorisation is to perform a task called *matrix completion*. Assume for a moment that you have a problem where you have a number of observations corresponding to different instances, but that these observations are both noisy and incomplete (you don't have observations of all features for all instances). You might ask if you can use unsupervised learning to *impute* the missing observations, and (potentially) remove some of the noise — this is exactly what matrix completion is about.

Of course, to perform matrix completion, you do have to make some kind of assumption about the processes that produced the data you have. One assumption that you could make is that the process is controlled by a small number of *latent variables*¹. Perhaps the simplest latent variable model is one based on low-rank factorisation: the chosen rank in such a case literally controls the number of latent variables, and thus the complexity of the model. We have known about such models since at least the late 1980s (where a technique called Latent Semantic Indexing was invented), but it wasn't until 2006² that we saw a gradient-based approach being proposed and used with a very large dataset.

The only difference between the pseudocode in Section 1, and one that works with incomplete observations is that we would only perform gradient updates where we have observations:

```
For epoch = 1 to N
  For r = 1 to m
    For c = 1 to n
      if  $\mathbf{A}_{rc}$  was observed
         $e = \mathbf{A}_{rc} - \hat{\mathbf{U}}_{r,*}(\hat{\mathbf{V}}_{c,*})^\top$ 
         $\hat{\mathbf{U}}_{r,*} \leftarrow \hat{\mathbf{U}}_{r,*} + \eta \cdot e \cdot \hat{\mathbf{V}}_{c,*}$ 
         $\hat{\mathbf{V}}_{c,*} \leftarrow \hat{\mathbf{V}}_{c,*} + \eta \cdot e \cdot \hat{\mathbf{U}}_{r,*}$ 
  return  $\hat{\mathbf{U}}, \hat{\mathbf{V}}$ 
```

¹you'll see many forms of latent variable model throughout this module, and also in advanced machine learning and in data-mining if you're taking them.

²<http://sifter.org/~simon/journal/20061211.html>

3.1 Implement and test masked factorisation (1 mark)

Create a new version of your factorisation code that additionally takes a second matrix of the same shape of the matrix being factorised. This second matrix should be a binary mask, with 1's where the value is valid, and zeros elsewhere. Use the mask to only compute updates using the valid values.

Your implementation should have the following signature:

```
def sgd_factorise_masked(A: torch.Tensor, M: torch.Tensor, rank: int,
                        num_epochs=1000, lr=0.01) -> Tuple[torch.Tensor, torch.Tensor]:
```

Use your function to compute the rank-2 factorisation of the following matrix (again using default learning rate and number of epochs):

$$\begin{bmatrix} 0.3374 & 0.6005 & 0.1735 \\ & 0.0492 & 1.8374 \\ 2.9407 & & 2.2620 \end{bmatrix}$$

What is your estimate of the completed matrix? How does this compare to the original matrix A ? What does this tell you?

4 Movie Recommendation

Matrix completion can be used to solve real-world problems, such as building a recommender system. In fact, the approach you implemented in the previous section arose because of a competition to build a recommender system for *Netflix*. The original *Netflix Challenge* provided ratings of movies (on a scale of 1 to 5) given by different viewers, and the challenge was to predict how a user might rate a film that they hadn't seen before. This challenge is a specific form of recommendation called *collaborative filtering*, which it turns out can be viewed as a matrix completion problem. Collaborative filtering approaches work on the basis that people often share similar tastes — the underlying assumption being that there are a small number of latent groupings of people with similar likes and dislikes.

Consider that the raw data consists of tuples (user, movie, rating). This data could be converted into a matrix where each row corresponds to a movie, each column to a user, and the elements to the rating given by a particular user to a particular movie. Clearly not all users will have watched and rated all movies, so there will be many missing entries. Matrix completion provides a way of estimating what values those missing entries should take — essentially providing guesses for how a user would rate a movie that they haven't yet watched based on the ratings of others with similar patterns of taste.

The stochastic, one-element at-a-time, algorithm you wrote in part 3, can be used to perform matrix completion for very large datasets (it does in fact mirror the basic approach that was originally developed for the Netflix Challenge), however, the python implementation is very very slow. The following code implements a vectorised gradient descent that updates all the values at once (by summing all the updates to the U and V matrices) and is faster:

```
def gd_factorise_masked(A: torch.Tensor, M: torch.Tensor, rank: int,
                       num_epochs: int=1000, lr: float=1e-5) -> Tuple[torch.Tensor, torch.Tensor]:
    U = torch.rand(A.shape[0], rank)
    V = torch.rand(A.shape[1], rank)

    for e in range(num_epochs):
        err = (A - U @ V.t()) * M
        U += lr * err @ V
        V += lr * err.t() @ U

    return U, V
```

4.1 Predict a movie rating (1 mark)

I've created a pytorch tensor that contains the ratings of the top-1000 movies from the top 10000 most active users within the original Netflix dataset. You can download a serialised version of this tensor from <http://ecs-vlc.github.io/COMP6258/labs/lab1/ratings.pt> and load it using `torch.load('ratings.pt')`. I've also created a list that maps the row-index of the ratings tensor to the title of each movie, and provided that as a CSV file at <http://ecs-vlc.github.io/COMP6258/labs/lab1/titles.csv>.

Using the `gd_factorise_masked` function (default learning rate and epochs, and a rank of 5), what rating do you predict the 5th user (column index 4) might give to the movie "A Beautiful Mind"? What is the *sum of squared errors* of the resultant matrix, computed over the valid (e.g. unmasked) values?