

the university of edinburgh

Applied Machine Learning (AML)

Recommender Systems

Oisin Mac Aodha • Siddharth N.

Recommender systems are systems which attempt to recommend items to users.
 e.g. movies, books, online ads, restaurants, etc.



- Recommender systems are systems which attempt to recommend items to users.
 e.g. movies, books, online ads, restaurants, etc.
- This is achieved by using information users provide about other items.
 - e.g. their past viewing/ purchasing behaviour, which movies they rated high or low, which ads they clicked on, etc.



- Recommender systems are systems which attempt to recommend items to users.
 e.g. movies, books, online ads, restaurants, etc.
- This is achieved by using information users provide about other items.
 - e.g. their past viewing/ purchasing behaviour, which movies they rated high or low, which ads they clicked on, etc.
- They are commonly found on websites or services that attempt to personalise content for users.



The Recommendation Task

• The central task we would like to solve is how to predict a user's rating (i.e. score) for the items they have *not* yet seen.



The Recommendation Task

- The central task we would like to solve is how to predict a user's rating (i.e. score) for the items they have *not* yet seen.
- We assume we have access to rating data from other users, but this is likely to be very sparse, i.e. the average user only rates a very small number of items.



Sources of Rating Data

Explicit Feedback

- Ask users to rate items, e.g. rating from 1 to 5, or like (+1) versus dislike (-1), etc.
- Can be hard to get this information in practice.



Sources of Rating Data

Explicit Feedback

- Ask users to rate items, e.g. rating from 1 to 5, or like (+1) versus dislike (-1), etc.
- Can be hard to get this information in practice.

Implicit Feedback

- Extracted from user actions, e.g. click on a video, watch until the end, ...
- Weaker form of supervision.



Sources of Rating Data

Explicit Feedback

- Ask users to rate items, e.g. rating from 1 to 5, or like (+1) versus dislike (-1), etc.
- Can be hard to get this information in practice.

Implicit Feedback

- Extracted from user actions, e.g. click on a video, watch until the end, ...
- Weaker form of supervision.

Assumption

We will assume that we have explicit feedback, and that any missing ratings are **missing** at random.



• Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}$.



• Suppose we observed data from four different users: $\mathcal{U} = \{u1, u2, u3, u4\}.$

users u1 🔘 u2 🔘 u3 🔘 u4 🔘



- Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}.$
- Who provided ratings for **five** different items: $I = \{i1, i2, i3, i4, i5\}$.





- Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}.$
- Who provided ratings for **five** different items: $I = \{i1, i2, i3, i4, i5\}$.
- Each rating can be one of five possible values: {1, 2, 3, 4, 5}.

lisers	items
	i 1
ul 🔵	i 2
u2 🔵	i 3
u3 🔵	
u4 🔵	
	0 15



- Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}.$
- Who provided ratings for **five** different items: $I = \{i1, i2, i3, i4, i5\}$.
- Each rating can be one of five possible values: {1, 2, 3, 4, 5}.





- Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}.$
- Who provided ratings for **five** different items: $I = \{i1, i2, i3, i4, i5\}$.
- Each rating can be one of five possible values: {1, 2, 3, 4, 5}.





- Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}.$
- Who provided ratings for **five** different items: $I = \{i1, i2, i3, i4, i5\}$.
- Each rating can be one of five possible values: {1, 2, 3, 4, 5}.





- Suppose we observed data from four different users: $U = \{u1, u2, u3, u4\}$.
- Who provided ratings for **five** different items: $I = \{i1, i2, i3, i4, i5\}$.
- Each rating can be one of five possible values: {1, 2, 3, 4, 5}.





• We can represent user ratings this as a matrix **Y** of size $|\mathcal{U}| \times |\mathcal{I}|$, i.e. num_users × num_items.







Missing Data

• Our ratings matrix can have many missing entries.



- Most users do not rate many items.
- As a result the data we observe can be very **sparse**.



• How can we predict missing ratings?





• How can we predict missing ratings?



- Suppose we have a user u4 who has provided the following ratings:
 - i1: 1 star
 - i3: 2 stars
 - i4: 4 stars



• How can we predict missing ratings?



- Suppose we have a user u4 who has provided the following ratings:
 i1: 1 star
 i3: 2 stars
 i4: 4 stars
- We would like to predict how they would rate items 2 and 5.



• How can we predict missing ratings?



- Suppose we have a user u4 who has provided the following ratings:
 i1: 1 star
 i3: 2 stars
 i4: 4 stars
- We would like to predict how they would rate items 2 and 5.
- We refer to the estimated ratings as \hat{Y}_{ui} .



• One simple approach is to report the **average** rating for an item.





• One simple approach is to report the **average** rating for an item.



• Average rating:

$$\hat{Y}_{ui} = k \sum_{\substack{u'=1\\Y_{u'i}\neq?}}^{|\mathcal{U}|} Y_{u'i}$$



• One simple approach is to report the **average** rating for an item.



• Average rating:

$$\hat{Y}_{ui} = k \sum_{\substack{u'=1\\Y_{u'i}\neq?}}^{|\mathcal{U}|} Y_{u'i}$$

• *k* is a normalisation factor:

$$k = 1 / \sum_{\substack{u'=1\\Y_{u'i} \neq ?}}^{|\mathcal{U}|} 1$$



• One simple approach is to report the **average** rating for an item.



• Average rating:

$$\hat{Y}_{ui} = k \sum_{\substack{u'=1\\Y_{u'i}\neq?}}^{|\mathcal{U}|} Y_{u'i}$$

• *k* is a normalisation factor:

$$k = 1 / \sum_{\substack{u'=1\\Y_{u'} \neq ?}}^{|\mathcal{U}|} 1$$

• This approach is simple, but it fails to capture differences between users.



Collaborative Filtering

- Another approach is to use rating information from **other** users.
- We can find **similar** users and then make predictions for our user of interest based their similarity to these existing users.



Collaborative Filtering

- Another approach is to use rating information from **other** users.
- We can find **similar** users and then make predictions for our user of interest based their similarity to these existing users.

Assumption

• The underlying assumption is that if user A and B rate items they have both seen similarly, then user A is more likely to rate items they have not seen similar to how user B would.



• We can **weight** the ratings score based on the **similarity** between the users.





• We can **weight** the ratings score based on the **similarity** between the users.



• Weighted average:





• We can **weight** the ratings score based on the **similarity** between the users.



• Weighted average:



• Here sim() is a similarity measure between a pair of users.



• We can **weight** the ratings score based on the **similarity** between the users.



• Weighted average:



- Here sim() is a similarity measure between a pair of users.
- *k* is a normalisation factor:

$$k = 1 / \sum_{u'=1}^{|\mathcal{U}|} \operatorname{sim}(u, u')$$



Limitations

- We need to define an appropriate similarity measure.
- Computing reliable similarity can be difficult with very sparse data.
- We need to store the entire dataset in memory at inference time.



Matrix Factorisation
Matrix Factorisation

- The recommendation task can be viewed as one of **matrix completion**.
- Here the goal is to predict all the missing entries of Y given the subset of user rating pairs (u, i) ∈ S we have observed (i.e. Y_{ui} ≠?)



Matrix Factorisation

- The recommendation task can be viewed as one of **matrix completion**.
- Here the goal is to predict all the missing entries of Y given the subset of user rating pairs (u, i) ∈ S we have observed (i.e. Y_{ui} ≠?)

$$\mathcal{L}(\hat{\mathbf{Y}}) = \sum_{(u,i)\in\mathcal{S}} (Y_{ui} - \hat{Y}_{ui})^2$$
$$= ||\mathbf{Y} - \hat{\mathbf{Y}}||^2$$



• The previous problem is under specified i.e. there are infinitely many ways of filing in the missing entries of *Y*.



- The previous problem is under specified i.e. there are infinitely many ways of filing in the missing entries of *Y*.
- We need to add some constraints.



- The previous problem is under specified i.e. there are infinitely many ways of filing in the missing entries of *Y*.
- We need to add some constraints.

Assumption

• Assume that **Y** is a low rank matrix.



- The previous problem is under specified i.e. there are infinitely many ways of filing in the missing entries of *Y*.
- We need to add some constraints.

Assumption

- Assume that **Y** is a low rank matrix.
- We can rewrite it as $\hat{Y} = UV^{\mathsf{T}} \approx Y$.



- The previous problem is under specified i.e. there are infinitely many ways of filing in the missing entries of *Y*.
- We need to add some constraints.

Assumption

- Assume that **Y** is a low rank matrix.
- We can rewrite it as $\hat{Y} = UV^{\mathsf{T}} \approx Y$.
- Here U is a $|\mathcal{U}| \times K$ matrix and V is a $|\mathcal{I}| \times K$ matrix.



Matrix Factorization

• Here we assume that our data can be represented as the product of two matrices $\hat{Y} = UV^{T} \approx Y$.





Matrix Factorization

• Here we assume that our data can be represented as the product of two matrices $\hat{Y} = UV^{T} \approx Y$.





Matrix Factorization

• Here we assume that our data can be represented as the product of two matrices $\hat{Y} = UV^{T} \approx Y$.





 We assume that our data can be represented as the product of a user matrix U and an item matrix V, i.e. Ŷ = UV^T ≈ Y.



- We assume that our data can be represented as the product of a user matrix U and an item matrix V, i.e. Ŷ = UV^T ≈ Y.
- Each row u_u of U represents a different user.



- We assume that our data can be represented as the product of a user matrix U and an item matrix V, i.e. Ŷ = UV^T ≈ Y.
- Each row u_u of U represents a different user.
- Similarly, each row v_i of V represents a different item.



- We assume that our data can be represented as the product of a user matrix U and an item matrix V, i.e. Ŷ = UV^T ≈ Y.
- Each row u_u of U represents a different **user**.
- Similarly, each row v_i of V represents a different item.
- To predict a missing entry we simply evaluate $\hat{Y}_{ui} = \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i$.



We can estimate the model weights
 θ = { U, V} using stochastic gradient descent (SGD).



- We can estimate the model weights
 θ = { U, V} using stochastic gradient descent (SGD).
- As there is missing data, we only have a subset of user rating pairs (u, i) ∈ S, where Y_{ui} ≠?.



- We can estimate the model weights
 θ = { U, V} using stochastic gradient descent (SGD).
- As there is missing data, we only have a subset of user rating pairs (u, i) ∈ S, where Y_{ui} ≠?.
- Our loss function is

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{(u,i)\in\mathcal{S}} (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)^2.$$



- We can estimate the model weights
 θ = { U, V} using stochastic gradient descent (SGD).
- As there is missing data, we only have a subset of user rating pairs (u, i) ∈ S, where Y_{ui} ≠?.
- Our loss function is

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{(u,i)\in\mathcal{S}} (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)^2.$$

Require: step size η , number steps N, valid indices S

- 1: $U \leftarrow initialisation$
- 2: $V \leftarrow$ initialisation
- 3: for $n \leftarrow 1$ to N do

4:
$$(u, i) \in \mathcal{S}$$

5:
$$e_{ui} = (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)$$

6:
$$\boldsymbol{u}_u \leftarrow \boldsymbol{u}_u + 2\eta (e_{ui} \boldsymbol{v}_i)$$

7:
$$v_i \leftarrow v_i + 2\eta (e_{ui}u_u)$$

8: return U, V



Regularisation

• We can also regularise our weights so that they do not become too large.



Regularisation

- We can also regularise our weights so that they do not become too large.
- Our loss function then becomes

$$\sum_{(u,i)\in\mathcal{S}} (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)^2 + \lambda (\sum_u ||\boldsymbol{u}_u||^2 + \sum_i ||\boldsymbol{v}_i||^2).$$



Regularisation

- We can also regularise our weights so that they do not become too large.
- Our loss function then becomes

$$\sum_{(u,i)\in\mathcal{S}} (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)^2 + \lambda (\sum_u ||\boldsymbol{u}_u||^2 + \sum_i ||\boldsymbol{v}_i||^2).$$

Require: step size η , number steps N, regularisation strength λ , valid indices S

- 1: $U \leftarrow$ initialisation
- 2: $V \leftarrow$ initialisation
- 3: for $n \leftarrow 1$ to N do

4:
$$(u, i) \in \mathcal{S}$$

5:
$$e_{ui} = (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)$$

6:
$$\boldsymbol{u}_u \leftarrow \boldsymbol{u}_u + 2\eta (e_{ui}\boldsymbol{v}_i - \lambda \boldsymbol{u}_u)$$

7:
$$v_i \leftarrow v_i + 2\eta (e_{ui}u_u - \lambda v_i)$$

8: return U, V



Example Overview

• Here we will walk through an example of applying matrix factorisation to predict unobserved user ratings for a small toy problem.



Example Overview

- Here we will walk through an example of applying matrix factorisation to predict unobserved user ratings for a small toy problem.
- We will use SGD to estimate the model weights $\theta = \{ U, V \}$.



Example Overview

- Here we will walk through an example of applying matrix factorisation to predict unobserved user ratings for a small toy problem.
- We will use SGD to estimate the model weights $\theta = \{ U, V \}$.
- Our training loss is the squared error with additional weight regularisation

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{(u,i)\in\mathcal{S}} (Y_{ui} - \boldsymbol{u}_u^{\mathsf{T}} \boldsymbol{v}_i)^2 + \lambda (\sum_u ||\boldsymbol{u}_u||^2 + \sum_i ||\boldsymbol{v}_i||^2).$$



• In this example we have 6 users and 5 items, where 13 of the ratings are missing.





- In this example we have 6 users and 5 items, where 13 of the ratings are missing.
- Our model will use 2 factors i.e. K = 2.





- In this example we have 6 users and 5 items, where 13 of the ratings are missing.
- Our model will use 2 factors i.e. K = 2.





Training Error

• Here we plot the error $|| \mathbf{Y} - \mathbf{U} \mathbf{V}^{\mathsf{T}} ||^2$ obtained during training.





• Below we see the learned factors after running SGD.





Matrix Factorisation - Example Results

• On the left we see the observed input data *Y*, and on the right we see the model predictions \hat{Y} .



items

c lach	4.4	4.8	4.6	1.0	-0.3
	4.9	4.5	5.5	3.1	1.9
	4.3	4.1	4.7	2.3	1.1
	2.4	0.8	3.5	5.1	4.9
	1.3	0.0	2.1	3.9	3.9
	1.1	-0.3	1.9	3.8	3.9

 $\hat{m{Y}}$



Visualising the Learned User Factors







Visualising the Learned Item Factors







The Netflix Prize

• In 2006 Netflix released a dataset containing 100,480,507 movie ratings that 480,189 users gave to 17,770 movies.



The Netflix Prize

- In 2006 Netflix released a dataset containing 100,480,507 movie ratings that 480,189 users gave to 17,770 movies.
- Each rating was an integer between 1 and 5, where a higher number indicated that a user liked a movie more.



The Netflix Prize

- In 2006 Netflix released a dataset containing 100,480,507 movie ratings that 480,189 users gave to 17,770 movies.
- Each rating was an integer between 1 and 5, where a higher number indicated that a user liked a movie more.
- They offered a prize of \$1,000,000 for the team that could improve prediction performance compared to the algorithm Netflix used at the time on a held out test set.


The Netflix Prize - Leaderboard

Netflix Prize					
Rule	es Leaderboard Register	Update	Submit	Download	
Lea	aderboard			Display top 4	0 leaders.
Rank	Team Name No Grand Prize candidates yet	Best	Score	% Improvement 	Last Submit Time
Grand	Prize - RMSE <= 0.8563				
1	PragmaticTheory	0.	8584	9.78	2009-06-16 01:04:47
	BellKor in BigChaos	0.	8590	9.71	2009-05-13 08:14:09
3	Grand Prize Team	0.	8593	9.68	2009-06-12 08:20:24
1	Dace	0.	8604	9.56	2009-04-22 05:57:03
5	BigChaos	0.	8613	9.47	2009-06-15 18:03:55
Progr	<u>ess Prize 2008</u> - RMSE = 0.86	16 - Winn	ing Tean	: BellKor in BigCl	iaos
	BellKor	0.	8620	9.40	2009-06-17 13:41:48
7	Gravity	0.	8634	9.25	2009-04-22 18:31:32
в	Opera Solutions	0.	8640	9.19	2009-06-09 22:24:53
9	xivector	0.	8640	9.19	2009-06-17 12:47:27

Image credit: https://www.wired.com/2012/04/netflix-prize-costs



The Netflix Prize - Data Privacy

• The original dataset did **not** provide any user identifying information, i.e. user names and IDs were anonymised.



The Netflix Prize - Data Privacy

- The original dataset did **not** provide any user identifying information, i.e. user names and IDs were anonymised.
- However, researchers were able to identify users in the dataset by matching their ratings against other publicly available sources of data (e.g. reviews they provided on the website IMDB).
- As a result, the Netflix Prize dataset is no longer publicly available.



Summary

• We discussed the recommendation problem.



Summary

- We discussed the recommendation problem.
- Recommendation systems are commonly deployed in services that aim to recommend items to users, e.g. movies, books, ads, ...



Summary

- We discussed the recommendation problem.
- Recommendation systems are commonly deployed in services that aim to recommend items to users, e.g. movies, books, ads, ...
- We outlined a memory-based collaborative filtering method and a model-based matrix factorisation approach for solving for unobserved ratings.

