

Problem 1 (30 points). Consider the following movie rating matrix with five users.

	LOTR	HPATPOZ	Snatch	LSATSB	The Gentlemen	The Hobbit
A	5	?	1	2	3	4
B	5	4	2	2	2	5
C	1	2	4	?	4	3
D	?	2	3	5	?	?
E	?	3	5	4	5	1

- Compute the user-user similarity for all the 10 pairs of users and 15 pairs of movies using Pearson's similarity. (Ignore missing values when computing similarity, i.e., you only need to consider the *commonly rated* entries, which can be a lower-dimensional vector).
- Let $k = 3$. Fill the missing entries of the matrix above using k -NN user-user CF and Pearson's similarity. When predicting, use the following formula:

$$\hat{r}_{u,i} = \bar{r}_u + \frac{\sum_{j \in K_u} S(u, u_j)(r_{u_j, i} - \bar{r}_{u_j})}{\sum_{j \in K_u} |S(u, u_j)|},$$

where \bar{r}_u is the average rating of user u (on the items that they actually have rated) and K_u is the top neighbours of u who also rated item i . Note here we rank user similarities by the absolute values of $S(u, u_j)$ since $S(u, u_j)$ can be negative for pearson similarity. (If I always like things you hate, then your rating is also very useful to me.)

(1) Pearson's similarity $s_p(u, v) = \frac{(\bar{r}_u - \hat{r}_u) \cdot (\bar{r}_v - \hat{r}_v)}{\|\bar{r}_u - \hat{r}_u\| \|\bar{r}_v - \hat{r}_v\|}$

The ave here is from common entries. Yuhan said this is okay because I've done it before announcement

$\bar{r}_A = (5, 1, 2, 3, 4)$	$\bar{r}_B = (5, 2, 2, 2, 5)$	
$\hat{r}_A = (3, 3, 3, 3, 3)$	$\hat{r}_B = (2.2, 3.2, 3.2, 3.2, 3.2)$	
$\bar{r}_A - \hat{r}_A = (2, -2, -1, 0, 1)$	$\bar{r}_B - \hat{r}_B = (-1.8, -1.2, -1.2, -1.2, 1.8)$	

Calculate using python

```
a = np.array([5,1,2,3,4])
b = np.array([5,2,2,2,5])

mean_a = np.mean(a)
mean_b = np.mean(b)

print(mean_a, mean_b)

a = a - mean_a
b = b - mean_b

print(a,b)

c = np.dot(a,b)
d = np.linalg.norm(a) * np.linalg.norm(b)
print(c,d)

similarity = c / d
print(similarity)

3.0 3.2
[ 2. -2. -1.  0.  1.] [ 1.8 -1.2 -1.2 -1.2  1.8]
9.0 10.392304845413266
0.8660254037844385
```

$s_p(A, B) =$

$$\frac{(2, -2, -1, 0, 1) \cdot (-1.8, -1.2, -1.2, -1.2, 1.8)}{\|(2, -2, -1, 0, 1)\| \|(-1.8, -1.2, -1.2, -1.2, 1.8)\|}$$

$$= \frac{2 \times -1.8 + -2 \times -1.2 + -1 \times -1.2 + 0 \times -1.2 + 1 \times 1.8}{\sqrt{2^2 + (-2)^2 + (-1)^2 + 0^2 + 1^2} \times \sqrt{(-1.8)^2 + (-1.2)^2 + (-1.2)^2 + (-1.2)^2 + (1.8)^2}}$$

$$= 0.866$$

$$AC : \bar{r}_A = (5, 1, 3, 4) \quad \bar{r}_C = (1, 4, 4, 3)$$

$$\hat{r}_A = (3.25, 3.25, 3.25, 3.25) \quad \hat{r}_C = (3, 3, 3, 3)$$

$$\bar{r}_A - \hat{r}_A = (1.75, -2.25, -0.25, 0.75) \quad \bar{r}_C - \hat{r}_C = (-2, 1, 1, 0)$$

$$S_p(A,C) = \frac{(1.75, -2.25, -0.25, 0.75) \cdot (-2, 1, 1, 0)}{\|(1.75, -2.25, -0.25, 0.75)\| \|(-2, 1, 1, 0)\|}$$

$$= -0.828$$

$$AD : \bar{r}_A = (1, 2) \quad \bar{r}_D = (3, 5)$$

$$\hat{r}_A = (1.5, 1.5) \quad \hat{r}_D = (4, 4)$$

$$\bar{r}_A - \hat{r}_A = (-0.5, 0.5) \quad \bar{r}_D - \hat{r}_D = (-1, 1)$$

$$S_p(A,D) = \frac{(-0.5, 0.5) \cdot (-1, 1)}{\|(-0.5, 0.5)\| \|(-1, 1)\|}$$

$$= 0.999$$

$$AE : \bar{r}_A = (1, 2, 3, 4) \quad \bar{r}_E = (5, 4, 5, 1)$$

$$\hat{r}_A = (2.5, 2.5, 2.5, 2.5) \quad \hat{r}_E = (3.75, 3.75, 3.75, 3.75)$$

$$\bar{r}_A - \hat{r}_A = (-1.5, -0.5, 0.5, 1.5) \quad \bar{r}_E - \hat{r}_E = (1.25, 0.25, 1.25, -2.75)$$

$$S_p(A,E) = \frac{(-1.5, -0.5, 0.5, 1.5) \cdot (1.25, 0.25, 1.25, -2.75)}{\|(-1.5, -0.5, 0.5, 1.5)\| \|(1.25, 0.25, 1.25, -2.75)\|}$$

$$= -0.750$$

$$BC : \bar{r}_B = (5, 4, 2, 2, 5) \quad \bar{r}_C = (1, 2, 4, 4, 3)$$

$$\hat{r}_B = (3.6, 3.6, 3.6, 3.6, 3.6) \quad \hat{r}_C = (2.8, 2.8, 2.8, 2.8, 2.8)$$

$$\bar{r}_B - \hat{r}_B = (1.4, 0.4, -1.6, -1.6, 1.4) \quad \bar{r}_C - \hat{r}_C = (-1.8, -0.8, 1.2, 1.2, 0.2)$$

$$S_P(\beta, c) = \frac{(-1.9, 0.4, -1.6, -1.6, 1.4) \cdot (-1.8, -0.8, 1.2, 1.2, 0.2)}{\|(-1.9, 0.4, -1.6, -1.6, 1.4)\| \|(-1.8, -0.8, 1.2, 1.2, 0.2)\|}$$

$$= \underline{-0.809}$$

$$\beta D: \bar{r}_B = (4, 2, 2) \quad \bar{r}_D = (2, 3, 5)$$

$$\hat{r}_B = (2.67, 2.67, 2.67) \quad \hat{r}_D = (1.33, 3.33, 3.33)$$

$$\bar{r}_D - \hat{r}_B = (1.33, -0.67, -0.67) \quad \bar{r}_D - \hat{r}_D = (-1.33, -0.33, 1.67)$$

$$S_P(\beta, D) = \frac{(1.33, -0.67, -0.67) \cdot (-1.33, -0.33, 1.67)}{\|(1.33, -0.67, -0.67)\| \|(-1.33, -0.33, 1.67)\|}$$

$$= \underline{-0.756}$$

$$\beta E: \bar{r}_B = (4, 2, 2, 2, 5) \quad \bar{r}_E = (3, 5, 4, 5, 1)$$

$$\hat{r}_B = (3, 3, 3, 3, 3) \quad \hat{r}_E = (3.1, 3.6, 3.6, 3.1, 3.6)$$

$$\bar{r}_E - \hat{r}_B = (1, -1, -1, -1, 2) \quad \bar{r}_E - \hat{r}_E = (-0.6, 1.4, 0.4, 1.4, -2.6)$$

$$S_P(\beta, E) = \frac{(1, -1, -1, -1, 2) \cdot (-0.6, 1.4, 0.4, 1.4, -2.6)}{\|(1, -1, -1, -1, 2)\| \|(-0.6, 1.4, 0.4, 1.4, -2.6)\|}$$

$$= \underline{-0.951}$$

$$CD: \bar{r}_C = (2, 4) \quad \bar{r}_D = (2, 3)$$

$$\hat{r}_C = (3, 3) \quad \hat{r}_D = (2.5, 2.5)$$

$$\bar{r}_D - \hat{r}_C = (-1, 1) \quad \bar{r}_D - \hat{r}_D = (-0.5, 0.5)$$

$$Sp(C, D) = \frac{(-1, 1) \cdot (-0.5, 0.5)}{\|(-1, 1)\| \|(-0.5, 0.5)\|}$$

$$= \underline{0.999}$$

$$CE : \bar{r}_C = (2, 4, 4, 3) \quad \bar{r}_E = (3, 5, 5, 1)$$

$$\hat{r}_C = (3.25, 3.25, 3.25, 3.25) \quad \hat{r}_E = (3.5, 3.5, 3.5, 3.5)$$

$$\bar{r}_C - \hat{r}_C = (-1.25, 0.75, 0.75, -0.25) \quad \bar{r}_E - \hat{r}_E = (-0.5, 1.5, 1.5, -2.5)$$

$$Sp(C, E) = \frac{(-1.25, 0.75, 0.75, -0.25) \cdot (-0.5, 1.5, 1.5, -2.5)}{\|(-1.25, 0.75, 0.75, -0.25)\| \|(-0.5, 1.5, 1.5, -2.5)\|}$$

$$= \underline{0.636}$$

$$DE : \bar{r}_D = (2, 3, 5) \quad \bar{r}_E = (3, 5, 4)$$

$$\hat{r}_D = (3.33, 3.33, 3.33) \quad \hat{r}_E = (4, 4, 4)$$

$$\bar{r}_D - \hat{r}_D = (-1.33, -0.33, 1.67) \quad \bar{r}_E - \hat{r}_E = (-1, 1, 0)$$

$$Sp(D, E) = \frac{(-1.33, -0.33, 1.67) \cdot (-1, 1, 0)}{\|(-1.33, -0.33, 1.67)\| \|(-1, 1, 0)\|}$$

$$= \underline{0.327}$$

Next, for simplicity, in calculating the Pearson's similarity for movie pairs, I will be showing the vector values, mean, and similarity.

We will represent movies with numbers for simplicity too. Calculated in python.

LOTR => 1

Snatch => 3

The gentleman => 5

HPI TPOZ => 2

LSATSB => 4

The Hobbit => 6

$$1-2 : \bar{r}_1 = (5, 1) \quad \bar{r}_2 = (4, 2)$$

$$\hat{r}_1 = 3 \quad \hat{r}_2 = 3$$

$$S_p(1,2) = 0.999$$

$$1-3 : \bar{r}_1 = (5, 5, 1) \quad \bar{r}_3 = (1, 2, 4)$$

$$\hat{r}_1 = 3.666 \quad \hat{r}_3 = 2.333$$

$$S_p(1,3) = -0.945$$

$$1-4 \quad \bar{r}_1 = (5, 5) \quad \bar{r}_4 = (2, 2)$$

$$\hat{r}_1 = 5 \quad \hat{r}_4 = 2$$

$$S_p(1,4) = 0$$

$$1-5 \quad \bar{r}_1 = (5, 5, 1) \quad \bar{r}_5 = (3, 2, 4)$$

$$\hat{r}_1 = 3.666 \quad \hat{r}_5 = 3$$

$$S_p(1,5) = -0.866$$

$$1-6 \quad \bar{r}_1 = (5, 5, 1) \quad \bar{r}_6 = (4, 5, 3)$$

$$\hat{r}_1 = 3.666 \quad \hat{r}_6 = 4$$

$$S_p(1,6) = 0.866$$

$$2-3 \quad \bar{r}_2 = (1, 2, 2, 3) \quad \bar{r}_3 = (2, 4, 3, 5)$$

$$\hat{r}_2 = 2.75 \quad \hat{r}_3 = 3.5$$

$$sp(2,3) = -0.405$$

$$2-4 \quad \bar{r}_2 = (1, 2, 3) \quad \bar{r}_4 = (2, 5, 4)$$

$$\hat{r}_2 = 3 \quad \hat{r}_4 = 3.666$$

$$sp(2,4) = -0.082$$

$$2-5 \quad \bar{r}_2 = (1, 2, 3) \quad \bar{r}_5 = (2, 4, 5)$$

$$\hat{r}_2 = 3 \quad \hat{r}_5 = 3.666$$

$$sp(2,5) = -0.654$$

$$2-6 \quad \bar{r}_2 = (1, 2, 3) \quad \bar{r}_6 = (5, 3, 1)$$

$$\hat{r}_2 = 3 \quad \hat{r}_6 = 3$$

$$sp(2,6) = 0.499$$

$$3-4 \quad \bar{r}_3 = (1, 2, 3, 5) \quad \bar{r}_4 = (2, 2, 5, 4)$$

$$\hat{r}_3 = 2.75 \quad \hat{r}_4 = 3.25$$

$$sp(3,4) = 0.683$$

$$3-5 \quad \bar{r}_3 = (1, 2, 4, 5) \quad \bar{r}_5 = (3, 2, 4, 5)$$

$$\hat{r}_3 = 3 \quad \hat{r}_5 = 3.5$$

$$sp(3,5) = 0.849$$

$$3-6 \quad \bar{r}_3 = (1, 2, 4, 5) \quad \bar{r}_6 = (4, 5, 3, 1)$$

$$\hat{r}_3 = 3 \quad \hat{r}_6 = 3.25$$

$$sp(3,6) = -0.855$$

$$4-5 \quad \bar{r}_4 = (2, 2, 4) \quad \bar{r}_5 = (3, 2, 5)$$

$$\hat{r}_4 = 2.667 \quad \hat{r}_5 = 3.333$$

$$sp(4,5) = 0.945$$

$$4-6 \quad \bar{r}_4 = (2, 2, 4) \quad \bar{r}_6 = (4, 5, 1)$$

$$\hat{r}_4 = 2.666 \quad \hat{r}_6 = 3.333$$

$$sp(4,6) = 0.971$$

$$5-6 \quad \bar{r}_5 = (3, 2, 4, 5) \quad \bar{r}_6 = (4, 5, 3, 1)$$

$$\hat{r}_5 = 3.5 \quad \hat{r}_6 = 3.25$$

$$sp(5,6) = -0.983$$

The ave used here is from Common entries.
 Yuhan said this is okay because
 I've done it before announcement, and it
 actually makes more sense. I am putting a note here that I am
 using a different mean as instructed by Yuhan.

(2)

$$\hat{r}_{u,i} = \bar{r}_u + \frac{\sum_{j \in K_u} S(u, u_j)(r_{u_j, i} - \bar{r}_{u_j})}{\sum_{j \in K_u} |S(u, u_j)|},$$

← in this formula, \hat{r} is the grade
and \bar{r} is average.

$$|L=3$$

$$A, 2 : \text{closest} \Rightarrow k \in \{B, C, D\}$$

$$A-B = 0.866, A-C = -0.828, A-D = 0.999$$

$$\bar{r}_{A,2} = 3 + \frac{0.866 \times (4 - 3.33) + (-0.828)(2 - 2.8) + (0.999)(2 - 3.33)}{|0.866| + |-0.828| + |0.999|}$$

$$= \underline{\underline{2.966}}$$

In this section, the mean is from all entries.
 But the "similarity" value uses mean from commonly rated entry.
 You have verified that this is okay as long as I'm consistent.

$$C, 4 : \text{closest} \Rightarrow k \in \{A, B, D\}$$

$$(-A) = -0.828$$

$$(-B) = -0.809$$

$$(-D) = 0.999$$

You have verified that this is okay
 as long as I'm consistent.

$$\bar{r}_{C,4} = 2.8 + \frac{-0.828(2 - 3) + (-0.809)(2 - 3.33) + 0.999(5 - 3.33)}{|-0.828| + |-0.809| + |0.999|}$$

$$= \underline{\underline{4.155}}$$

$$D, 1 : \text{closest } 3 \Rightarrow k \in \{A, B, C\}$$

$$(-A) = 0.999$$

$$(-B) = -0.756$$

$$(-C) = 0.999$$

$$\bar{r}_{D,1} = 3.333 + \frac{0.999(5 - 3) + (-0.756)(5 - 3.33) + 0.999(1 - 2.8)}{|0.999| + |-0.756| + |0.999|}$$

$$= \underline{\underline{2.947}}$$

$$D, 5 : \{ \text{closed} \} \Rightarrow k \in \{ A, B, D \}$$

$$D-A = 0.999$$

$$D-B = -0.756$$

$$D-C = 0.999$$

$$\bar{F}_{D,5} = 3.333 + \frac{0.999(3-3) - 0.756(2-3.33) + 0.999(4-2.8)}{|0.999| + |-0.756| + |0.999|}$$

$$= \underline{\underline{9.133}}$$

$$D, 6 : \{ \text{closed} \} \Rightarrow k \in \{ A, B, C \}$$

$$D-B = -0.756$$

$$D-C = 0.999$$

$$D-A = 0.999$$

$$\bar{F}_{D,6} = 3.333 + \frac{0.999(4-3) - 0.756(5-3.33) + 0.999(3-2.8)}{|0.999| + |-0.756| + |0.999|}$$

$$= \underline{\underline{3.311}}$$

$$E, 1 : \{ \text{closed} \} \Rightarrow k \in \{ A, B, C \}$$

$$E-A = -0.750$$

$$E-B = -0.951$$

$$E-C = 0.636$$

$$\bar{F}_{E,1} = 3.6 + \frac{(-0.750)(5-3) - 0.951(5-3.33) + 0.636(1-2.8)}{|-0.750| + |-0.951| + |0.636|}$$

$$= \underline{\underline{1.790}}$$

The answer may be different from the answer key because the Similarity calculated in 1.1 is using mean from commonly rated entries. However, the mean used for the formula in 1.2 uses the mean of all entries. So there are 2 types of mean involved. Please check my working instead of final answer. Thank you.

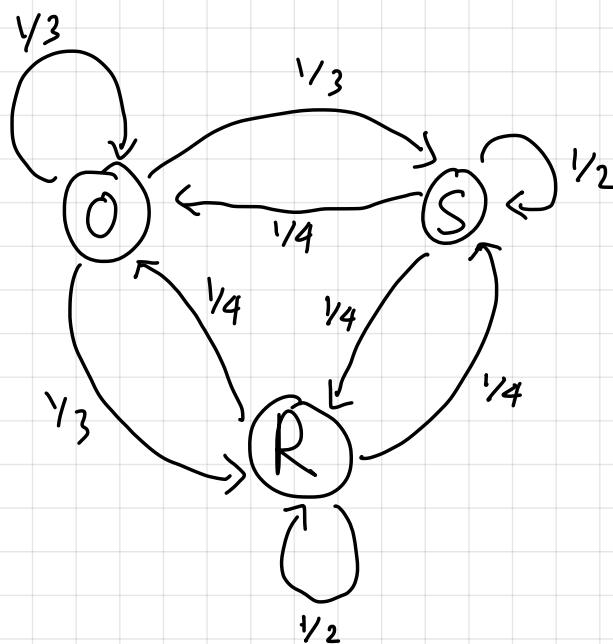
Problem 2 (30 points). Consider a Markov chain with three states, Overcast, Rain, and Sunny. The transition probabilities are given in the following table. The (i, j) th entry of the matrix is the probability that the next day to be j if today is i . November 29, 2020 is Rain.

	1	2	3
1	O	S	R
2	S		
3	R		

$$\begin{aligned} O &\rightarrow 1 \\ S &\rightarrow 2 \\ R &\rightarrow 3 \end{aligned}$$

1. Draw the state transition diagram with arrows annotating the transition probabilities.
2. What is the probability that it will be Sunny on November 30th, 2020?
3. What is the probability that it will Rain on December 2nd, 2020?
4. What is the probability that it will Rain every day until December 5, 2020 (including it)?
5. Compute the probability that it will Rain on December 6, 2020?

①



② $P(X_{t+1} = S | X_t = R) = \frac{1}{4}$ $\leftarrow P_{32}^{(1)} = \frac{1}{4}$

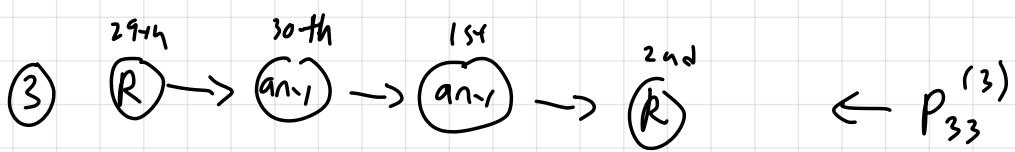
or

Find $P_{32}^{(1)}$

$$\begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/4 & 1/2 & 1/4 \\ 1/4 & 1/4 & 1/2 \end{bmatrix} = \begin{bmatrix} 1/4 & 1/4 & 1/2 \end{bmatrix} \Rightarrow P_{32}^{(1)} = \frac{1}{4}$$

↑
index 3
for R

↑
index 2
for S



$$P^3 = \begin{bmatrix} 0.27315 & 0.36343 & 0.36343 \\ 0.27257 & 0.37153 & 0.35590 \\ 0.27257 & 0.35590 & 0.37153 \end{bmatrix} \leftarrow \text{calculated in python}$$

```
import numpy as np
from numpy.linalg import matrix_power

temp = np.array([[1/3, 1/3, 1/3], [1/4, 1/2, 1/4], [1/4, 1/4, 1/2]])
print(temp)
power3 = matrix_power(temp, 3)
print(power3)

[[0.33333333 0.33333333 0.33333333]
 [0.25 0.5 0.25]
 [0.25 0.25 0.5]]
 [[0.27314815 0.36342593 0.36342593]
 [0.27256944 0.37152778 0.35590278]
 [0.27256944 0.35590278 0.37152778]]
```

$$P_{33}^{(3)} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.27315 & 0.36343 & 0.36343 \\ 0.27257 & 0.37153 & 0.35590 \\ 0.27257 & 0.35590 & 0.37153 \end{bmatrix}$$

$$= \begin{bmatrix} 0.27257 & 0.35590 & \cancel{0.37153} \\ \cancel{0.37153} & \cancel{0.37153} & \cancel{0.37153} \end{bmatrix} \Rightarrow \underline{\underline{0.37153}}$$

$\xrightarrow{3rd \text{ entry}}$

(4) $P(x_{t+6} = R | x_{t+5} = R) \times P(x_{t+5} = R | x_{t+4} = R) \times P(x_{t+4} = R | x_{t+3} = R)$

$$P(x_{t+3} = R | x_{t+2} = R) \times P(x_{t+2} = R | x_{t+1} = R) \times P(x_{t+1} = R | x_t = R)$$

$$= \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \underline{\underline{\frac{1}{64}}}$$

$$⑤ P_{33}^{(7)} = \text{3rd index of } \underbrace{U \cdot P^7}_{\substack{\text{initial state}}} \quad \leftarrow \text{calculated in Python}$$

$$P^7 = \begin{bmatrix} 0.27273 & 0.36364 & 0.36364 \\ 0.27273 & 0.36367 & 0.36361 \\ 0.27273 & 0.36361 & 0.36367 \end{bmatrix}$$

```
temp = np.array([[1/3, 1/3, 1/3], [1/4, 1/2, 1/4], [1/4, 1/4, 1/2]])
print(temp)
power7 = matrix_power(temp, 7)
print(power7)

[[0.33333333 0.33333333 0.33333333]
 [0.25 0.5 0.25]
 [0.25 0.25 0.5]]
 [[0.27272729 0.36363635 0.36363635]
 [0.27272727 0.36366689 0.36360585]
 [0.27272727 0.36360585 0.36366689]]
```

$$P_{33}^{(7)} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \cdot P^7 = \begin{bmatrix} 0.27273 & 0.36361 & \underbrace{0.36367}_{\substack{\dots}} \end{bmatrix}$$

3rd entry

$$\Rightarrow \underline{\underline{0.36367}}$$