

Direction Régionale : Casablanca – Settat Cité des métiers et des compétences Nouaceur

Contrôle continu N° 1 : Algorithmes distribués (M204)
Année 2025/2026

Filière : Intelligence Artificielle
Niveau : TS
Groupe : DIA_IAOBD_TS_201
Nom et Prénom :.....

Durée : 2H
Variante : A
Barème : .../20

Partie Pratique : (20 points)

Données doivent être disponibles dans HDFS sur le dossier : input_Nom_Stagiaire/movielens
movies.dat, ratings.dat, tags.dat

Structure du Dataset

- **movies.dat** : MovieID::Title::Genres
1 1::Toy Story (1995)::Animation|Children's|
2
 - **ratings.dat** : UserID::MovieID::Rating::Timestamp
1 1::122::5.0::838985046
2
 - **tags.dat** : UserID::MovieID::Tag::Timestamp
1 15::4973::excellent!::1215184630
2

1 Exercice 1 : MapReduce avec Hadoop (8 points)

Développer des programmes MapReduce en Python pour analyser le dataset MovieLens. Tous les scripts doivent être exécutés via Hadoop Streaming.

Q1.1 - Comptage des films par année (1 point)

Compter le nombre de films sortis chaque année.

- a) Écrire `mapper_year.py` qui émet (`année, 1`)
 - b) Écrire `reducer_year.py` qui agrège les résultats
 - c) Exécuter avec Hadoop Streaming

Sortie attendue : Liste année

nombre_films

```

1 #!/usr/bin/env python
2 import sys
3 import re
4 for line in sys.stdin:
5     line = line.strip()
6     fields = line.split('::')
7
8 -----
9
10 -----
11
12 -----
13
14 -----
15
16 -----
17
18 -----
19
20 -----
21
22 -----
23
24 -----

```

Code à compléter - mapper_year.py :

```

1 #!/usr/bin/env python
2 import sys
3
4 current_year = None
5 count = 0
6 for line in sys.stdin:
7     line = line.strip()
8     year, value = line.split('\t')
9
10 -----
11
12 -----
13
14 -----
15
16 -----
17
18 -----
19
20 -----
21
22 -----
23
24 -----

```

Code à compléter - reducer_year.py :

```

1 hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar \
2   -input /user/cloudera/movielens/movies.dat \
3   -output _____ \
4   -mapper _____ \
5   -reducer _____ \
6   -file mapper_year.py \
7   -file reducer_year.py

```

Commande d'exécution :

Q1.2 - Nombre d'évaluations par utilisateur (1,5 points)

Calculer le nombre d'évaluations faites par chaque utilisateur.

- Écrire `mapper_user_ratings.py` qui lit `ratings.dat`
- Émettre (`UserID`, 1) pour chaque évaluation
- Écrire `reducer_user_ratings.py` qui fait la somme
- Trier les résultats par nombre décroissant

Sortie attendue : Top 10 des utilisateurs les plus actifs

```
1 #!/usr/bin/env python
2 import sys
3
4 for line in sys.stdin:
5     line = line.strip()
6     fields = line.split('::')
7
8     -----
9
10    -----
11
12    -----
13
14    -----
15
16    -----
17
```

Code à compléter - `mapper_user_ratings.py` :

Q1.3 - Note moyenne par film (2 points)

Calculer la note moyenne de chaque film.

- Écrire `mapper_movie_avg.py` qui émet (`MovieID`, (rating, 1))
- Écrire `reducer_movie_avg.py` qui calcule la moyenne
- Filtrer les films avec au moins 100 évaluations

Indice : Utiliser un format de sortie `MovieID` `note_moyenne` `nombre_notes`

```
1 #!/usr/bin/env python
2 import sys
3
4 for line in sys.stdin:
5     line = line.strip()
6     fields = line.split('::')
7
8     -----
9
10    -----
11
12    -----
13
14    -----
15
16    -----
17
```

Code à compléter - `mapper_movie_avg.py` :

Q1.4 - Films par genre (1,5 points)

Compter le nombre de films pour chaque genre.

- Écrire `mapper_genre.py` qui éclate les genres multiples
- Pour un film avec genres "Action|Comedy|Drama", émettre 3 paires
- Écrire `reducer_genre.py` qui agrège les résultats

Sortie attendue : Liste des genres avec leur nombre de films

```
1 #!/usr/bin/env python
2 import sys
3
4 for line in sys.stdin:
5     line = line.strip()
6     fields = line.split('::')
7
8     -----
9
10    -----
11
12    -----
13
14    -----
15
16    -----
17
18    -----
19
20    -----
```

Code à compléter - `mapper_genre.py` :

Q1.5 - Top 10 des meilleurs films (2 points)

Trouver les 10 films avec la meilleure note moyenne (min 500 évaluations).

- Utiliser le résultat de Q1.3
- Écrire un second MapReduce qui trie les films par note moyenne
- Filtrer les films avec plus de 500 évaluations
- Produire un top 10

Sortie attendue : MovieID Titre Note_moyenne Nombre_notes

```
1 #!/usr/bin/env python
2 import sys
3
4 for line in sys.stdin:
5     line = line.strip()
6     # Format: MovieID    avg_rating    count
7     movie_id, avg_rating, count = line.split('\t')
8
9     -----
10    -----
11
12    -----
13
14    -----
15
16    -----
```

Code à compléter - `mapper_top10.py` :

2 Exercice 2 : Apache Pig (6 points)

Vous devez écrire des scripts Pig Latin pour analyser les données. Tous les résultats doivent être stockés dans HDFS.

Q2.1 - Chargement et inspection (1 point)

Charger les trois fichiers et afficher leur schéma.

Écrire un script Pig qui :

- Charge `movies.dat`, `ratings.dat`, `tags.dat` avec le bon séparateur
- Affiche le schéma de chaque relation
- Affiche 5 tuples de chaque

```
1 -- Chargement des films
2 -----
3 -----
4 -----
5 -----
6 -----
7 -----
8 -----
9 -----
10 -- Chargement des notes
11 -----
12 -----
13 -----
14 -----
15 -----
16 -----
17 -----
18 -- Chargement des tags
19 -----
20 -----
21 -----
22 -----
23 -----
24 -----
25 -----
26 -----
27 -- Afficher le schéma (DESCRIBE, DUMP, ...)
28 -----
29 -----
30 -----
31 -----
32 -----
33 -----
34 -----
```

Script Pig à compléter :

Q2.2 - Films les plus populaires (1 point)

Trouver les 20 films avec le plus d'évaluations.

Écrire un script Pig qui :

- Joint `movies` et `ratings` par `MovieID`
- Calcule le nombre d'évaluations par film
- Trie par nombre décroissant
- Stocke le résultat dans HDFS

```

1 -- Jointure films-notes
2 movies_ratings = -----
3
4
5 -- Regrouper par film
6 -----
7
8 -- Calculer le nombre d'évaluations
9 movies_count = FOREACH movies_group GENERATE
10   group AS movieId,
11   COUNT(movies_ratings) AS rating_count;
12
13 -- Trier et limiter (ORDER movies_count BY rating_count DESC) 20
14 -----
15
16 -----
17
18 -- Stocker dans HDFS
19 STORE top_movies INTO '-----';
20
21 -----

```

Script Pig à compléter :

Q2.3 - Analyse des tags (2 points)

Analyser la fréquence des tags.

Écrire un script Pig qui :

- Charge tags.dat
- Convertit tous les tags en minuscules
- Compte la fréquence de chaque tag
- Stocke les 50 tags les plus fréquents dans HDFS

```

1 -- Chargement des tags
2 tags = LOAD '-----';
3
4   USING PigStorage('::')
5
6   AS (userId:int, movieId:int, tag:chararray, timestamp:chararray);
7
8 -- Convertir en minuscules et compter (tags_lower)
9
10 -----
11
12 -----
13
14 -----
15
16 -----
17
18 -----
19
20 -- Top 50 (LIMIT (ORDER tags_count BY frequency DESC) 50)
21
22
23 -----
24
25 STORE -----

```

Script Pig à compléter :

Q2.4 - Utilisateurs les plus actifs par année (2 points)

Identifier les utilisateurs les plus actifs chaque année.

Écrire un script Pig qui :

- a) Extrait l'année du timestamp des évaluations
- b) Compte le nombre d'évaluations par utilisateur et par année
- c) Pour chaque année, garde le top 3 des utilisateurs
- d) Stocke le résultat dans HDFS

```
1 -- TODO Extraire l'année du timestamp
2 ratings_with_year = FOREACH ratings GENERATE
3 -----
4 -----
5 -----
6 -----
7 -----
8 -----
9 -----
10 -----
11 -----
12 -----
13 -----
14 -----
15 -----
16 -----
17 -----
18 -----
19 -----
20 -----
21 -----
22 -----
23 -----
24 -----
25 -----
26 -----
27 -----
28 -----
29 -----
30 -----
31 -----
32 -----
33 -----
```

Script Pig à compléter :

3 Exercice 3 : Apache Hive (6 points)

Vous devez créer des tables Hive et exécuter des requêtes analytiques.

Q3.1 - Création des tables externes (1 point)

Créer les tables Hive externes pointant vers les données HDFS.

Écrire les requêtes HiveQL pour :

- a) Créer la table `movies_hive` avec séparateur ::
- b) Créer la table `ratings_hive` avec séparateur ::
- c) Créer la table `tags_hive` avec séparateur ::
- d) Vérifier le nombre d'enregistrements dans chaque table

```
1 -- Table movies_hive  
2 -----  
3  
4 -----  
5  
6 -----  
7  
8 -----  
9  
10 -- Table ratings_hive  
11 -----  
12  
13 -----  
14  
15 -----  
16  
17 -----  
18  
19 -- Table tags_hive  
20 -----  
21  
22 -----  
23  
24 -----  
25  
26 -----  
27  
28 -- Vérification
```

Requête HiveQL à compléter :

Q3.2 - Analyse temporelle (1 point)

Analyser l'évolution des évaluations dans le temps.

Écrire une requête HiveQL qui :

- Calcule le nombre d'évaluations par mois (année-mois)
- Calcule la note moyenne par mois
- Affiche les résultats triés par période
- Sauvegarde le résultat dans une table Hive **ratings_monthly_stats**

```
1 -----  
2  
3 -----  
4  
5 -----  
6  
7 -----  
8  
9 -----  
10  
11 -----  
12  
13 -----  
14  
15 -----  
16  
17 -----
```

Requête HiveQL à compléter :

Q3.3 - Segmentation des utilisateurs (1 point)

Segmenter les utilisateurs selon leur comportement.

Écrire une requête HiveQL qui :

- Classe les utilisateurs en 4 catégories
- Pour chaque segment, calcule les statistiques
- Sauvegarde dans une table Hive `user_segments`

```
1 CREATE TABLE user_segments AS
2 WITH user_stats AS (
3     SELECT
4         userId,
5         COUNT(*) AS rating_count,
6         ROUND(AVG(rating), 3) AS avg_rating
7     FROM ratings_hive
8     GROUP BY userId
9 )
10 SELECT
11     CASE
12         WHEN rating_count >= 500 THEN 'Power Users'
13         WHEN rating_count >= 100 THEN 'Active Users'
14         WHEN rating_count >= 10 THEN 'Regular Users'
15         ELSE 'Casual Users'
16     END AS user_segment,
17     COUNT(*) AS user_count,
18     ROUND(AVG(rating_count), 1) AS avg_ratings_per_user,
19     ROUND(AVG(avg_rating), 3) AS segment_avg_rating
20 FROM user_stats
21
22 GROUP BY -----
23
24 ORDER BY ----- DESC;
25
26 -----
27
28 -----
```

Requête HiveQL à compléter:

Q3.4 - Analyse des genres (1 point)

Analyser la performance des différents genres.

Écrire une requête HiveQL qui :

- Éclate les genres multiples
- Calcule pour chaque genre les statistiques
- Filtre les genres avec au moins 100 films
- Sauvegarde dans une table Hive `genre_performance`

```
1 CREATE TABLE genre_performance AS
2 SELECT
3     genre,
4     COUNT(DISTINCT m.movieId) AS film_count,
5     ROUND(AVG(r.rating), 3) AS avg_rating,
6     COUNT(r.rating) AS total_ratings
7
8 -----
9
10 -----
```

11
12
13
14
15
16
17
18
19
20

Requête HiveQL à compléter :

Q3.5 - Recommandation avancée (2 points)

Implémenter un système de recommandation collaborative simple.

Écrire une requête HiveQL qui pour un film donné (ex : MovieID=1) :

- a) Trouve les utilisateurs qui ont aimé ce film
- b) Trouve les autres films que ces utilisateurs ont aussi aimés
- c) Calcule un score de recommandation
- d) Recommande les 10 films avec le score le plus élevé

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Requête HiveQL à compléter :