

# Projet CinéExplorer

Livrable 1 : Exploration et Base SQLite

OUDELET Clément  
4A Info - Polytech Marseille

3 décembre 2025

## Table des matières

<b>1</b>	<b>Exploration des Données (T1.0)</b>	<b>2</b>
1.1	Observations Clés . . . . .	2
1.2	Visualisations . . . . .	2
<b>2</b>	<b>Conception du Schéma Relationnel (T1.1)</b>	<b>2</b>
<b>3</b>	<b>Requêtes SQL (T1.3)</b>	<b>3</b>
3.1	Q1 : Filmographie d'un acteur . . . . .	3
3.2	Q2 : Top N films d'un genre . . . . .	3
3.3	Q3 : Acteurs multi-rôles . . . . .	3
3.4	Q4 : Collaborations Réalisateur-Acteur . . . . .	3
3.5	Q5 : Genres populaires . . . . .	4
3.6	Q6 : Évolution de carrière (CTE) . . . . .	4
3.7	Q7 : Classement par genre (Window Function) . . . . .	4
3.8	Q8 : Carrière propulsée . . . . .	5
3.9	Q9 : Requête Libre (Complexe) . . . . .	5
<b>4</b>	<b>Indexation et Benchmark (T1.4)</b>	<b>6</b>
4.1	Stratégie d'Indexation . . . . .	6
4.2	Résultats du Benchmark . . . . .	6
<b>5</b>	<b>Conclusion Phase 1</b>	<b>6</b>

# 1 Exploration des Données (T1.0)

## 1.1 Observations Clés

L'analyse exploratoire réalisée avec Pandas a révélé plusieurs caractéristiques importantes du dataset :

- **Volumétrie** : Le dataset contient plusieurs milliers de films et d'acteurs. La table de liaison principaux est la plus dense.
- **Qualité des données** : Les fichiers utilisent le caractère \N pour désigner les valeurs nulles, traité lors de l'import.
- **Cohérence** : Des tests d'intégrité ont validé la correspondance entre les clés étrangères (mid, pid) à travers les fichiers.

## 1.2 Visualisations

Voici la distribution temporelle des films et les genres les plus représentés.

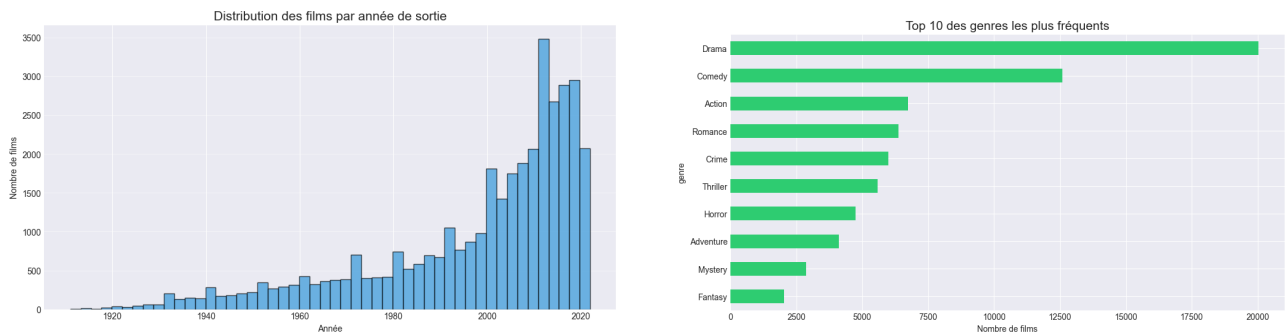


FIGURE 1 – Distribution des films par année et Top 10 des genres

# 2 Conception du Schéma Relationnel (T1.1)

Le schéma a été conçu en **3ème Forme Normale (3NF)**. Les listes séparées par des virgules (Genres, Titres alternatifs) ont été atomisées dans des tables dédiées.

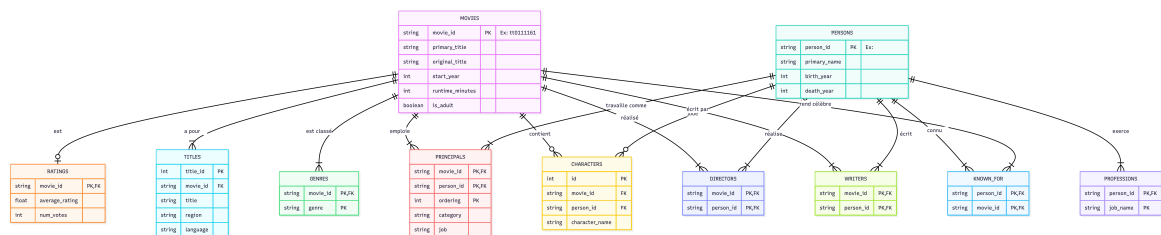


FIGURE 2 – Diagramme Entité-Relation (ER) du projet

## 3 Requêtes SQL (T1.3)

Nous avons implémenté les 9 requêtes demandées, couvrant des concepts SQL avancés.

### 3.1 Q1 : Filmographie d'un acteur

*Objectif : Récupérer films, rôles et notes pour un acteur (recherche souple).*

Récupération des données via jointures (JOIN) entre `persons`, `principals` et `movies`. L'utilisation de `LEFT JOIN` préserve les films même si les données annexes (rôle, note) sont manquantes, et `LIKE` permet une recherche approximative.

```
1 SELECT m.primary_title, m.start_year, c.character_name, r.average_rating
2 FROM movies m
3 JOIN principals p ON m.movie_id = p.movie_id
4 JOIN persons pe ON p.person_id = pe.person_id
5 LEFT JOIN characters c ON m.movie_id = c.movie_id AND pe.person_id = c.person_id
6 LEFT JOIN ratings r ON m.movie_id = r.movie_id
7 WHERE pe.primary_name LIKE ?
8 ORDER BY m.start_year DESC
```

### 3.2 Q2 : Top N films d'un genre

*Objectif : Filtrage multicritères (année, genre) et tri sur la note.*

Filtrage simultané sur trois tables optimisé par les index sur `genres` et `start_year`. Le tri descendant sur la note permet d'extraire efficacement le Top N.

```
1 SELECT m.primary_title, m.start_year, r.average_rating
2 FROM movies m
3 JOIN genres g ON m.movie_id = g.movie_id
4 JOIN ratings r ON m.movie_id = r.movie_id
5 WHERE g.genre = ? AND m.start_year BETWEEN ? AND ?
6 ORDER BY r.average_rating DESC LIMIT ?
```

### 3.3 Q3 : Acteurs multi-rôles

*Objectif : Identifier les performances multiples dans un même film (ex : Dr. Folamour).*

Détection des doublons sémantiques par agrégation `GROUP BY` (film, acteur) et filtrage post-agrégation avec `HAVING count > 1`.

```
1 SELECT pe.primary_name, m.primary_title, COUNT(c.character_name) as nb
2 FROM characters c
3 JOIN persons pe ON c.person_id = pe.person_id
4 JOIN movies m ON c.movie_id = m.movie_id
5 GROUP BY c.movie_id, c.person_id
6 HAVING nb > 1
7 ORDER BY nb DESC
```

### 3.4 Q4 : Collaborations Réalisateur-Acteur

*Objectif : Double jointure pour relier deux personnes via un film commun.*

Identification des paires via jointure entre la filmographie de l'acteur (sous-requête) et la table `directors` sur le champ commun `movie_id`.

```
1 SELECT d_pe.primary_name, COUNT(*) as collab
2 FROM principals p_actor
3 JOIN directors d ON p_actor.movie_id = d.movie_id
4 JOIN persons d_pe ON d.person_id = d_pe.person_id
5 WHERE p_actor.person_id = (SELECT person_id FROM persons WHERE primary_name LIKE ?)
6 GROUP BY d.person_id
7 ORDER BY collab DESC
```

### 3.5 Q5 : Genres populaires

*Objectif : Agrégation (Moyenne) avec condition sur le groupe (HAVING).*

Calcul de statistiques agrégées par genre. Le filtre **HAVING** garantit la pertinence statistique en excluant les genres ayant un échantillon trop faible (< 50 films).

```
1 SELECT g.genre, AVG(r.average_rating) as avg_rat, COUNT(*) as nb
2 FROM genres g
3 JOIN ratings r ON g.movie_id = r.movie_id
4 GROUP BY g.genre
5 HAVING avg_rat > 7.0 AND nb > 50
6 ORDER BY avg_rat DESC
```

### 3.6 Q6 : Évolution de carrière (CTE)

*Objectif : Utilisation d'une Common Table Expression pour grouper par décennie.*

Utilisation d'une **CTE** pour pré-calculer la décennie de chaque film, simplifiant l'agrégation finale et améliorant la lisibilité du code.

```
1 WITH career AS (
2     SELECT
3         (m.start_year / 10) * 10 as decade,
4         r.average_rating
5     FROM movies m
6     JOIN principals p ON m.movie_id = p.movie_id
7     JOIN persons pe ON p.person_id = pe.person_id
8     LEFT JOIN ratings r ON m.movie_id = r.movie_id
9     WHERE pe.primary_name = ? AND m.start_year IS NOT NULL
10 )
11 SELECT decade, COUNT(*) as count, AVG(average_rating)
12 FROM career
13 GROUP BY decade
14 ORDER BY decade
```

### 3.7 Q7 : Classement par genre (Window Function)

*Objectif : Classement local (par genre) sans sous-requête corrélée.*

Classement performant en une seule passe grâce à la fonction de fenêtrage **RANK()** partitionnée par genre, permettant de filtrer directement le Top 3.

```
1 WITH ranked AS (
2     SELECT
3         g.genre,
4         m.primary_title,
5         r.average_rating,
6         RANK() OVER (PARTITION BY g.genre ORDER BY r.average_rating DESC) as rk
7     FROM movies m
8     JOIN genres g ON m.movie_id = g.movie_id
9     JOIN ratings r ON m.movie_id = r.movie_id
10     -- WHERE r.num_votes > 5000
11 )
12 SELECT genre, primary_title, average_rating, rk
13 FROM ranked
14 WHERE rk <= 3
```

### 3.8 Q8 : Carrière propulsée

*Objectif : Identification de films à fort impact (>200k votes).*

Isolation du premier succès chronologique (>200k votes) pour chaque acteur via ROW\_NUMBER() ordonné par date, identifiant ainsi le rôle charnière.

```
1 WITH FirstHit AS (  
2     SELECT  
3         pe.primary_name,  
4         m.primary_title,  
5         m.start_year,  
6         r.num_votes,  
7         -- On classe par annee SEULEMENT les films > 200k pour chaque acteur  
8         ROW_NUMBER() OVER (  
9             PARTITION BY pe.person_id  
10            ORDER BY m.start_year ASC  
11        ) as hit_rank  
12 FROM persons pe  
13 JOIN principals p ON pe.person_id = p.person_id  
14 JOIN movies m ON p.movie_id = m.movie_id  
15 JOIN ratings r ON m.movie_id = r.movie_id  
16 WHERE r.num_votes > 200000  
17        AND p.category IN ('actor', 'actress') -- On cible les acteurs  
18 )  
19 SELECT primary_name, primary_title, start_year, num_votes  
20 FROM FirstHit  
21 WHERE hit_rank = 1 -- On prend uniquement leur PREMIER gros succes  
22 ORDER BY num_votes DESC  
23 LIMIT 20
```

### 3.9 Q9 : Requête Libre (Complexe)

*Objectif : Films longs (>2h), très bien notés (>8.0) et en langue française.*

Intersection stricte de trois critères (technique, critique, régional) réalisée via des jointures internes (INNER JOIN) sur les tables correspondantes.

```
1 SELECT m.primary_title, m.runtime_minutes, r.average_rating  
2 FROM movies m  
3 JOIN ratings r ON m.movie_id = r.movie_id  
4 JOIN titles t ON m.movie_id = t.movie_id  
5 WHERE m.runtime_minutes > 120  
6        AND r.average_rating > 8.0  
7        AND t.language = 'fr'  
8 LIMIT 20
```

## 4 Indexation et Benchmark (T1.4)

### 4.1 Stratégie d'Indexation

Pour optimiser les performances, nous avons créé des index sur les colonnes de filtrage (`primary_name`, `genre`, `start_year`) et sur toutes les clés étrangères pour accélérer les jointures.

### 4.2 Résultats du Benchmark

Les tests ont été effectués sur un MacBook Air (Puce M3 16go RAM). Les résultats montrent un gain massif sur les requêtes lourdes (Q4, Q6), tandis que l'optimiseur de requêtes SQLite sur les petits ensembles de données (Q3, Q8) montre parfois un léger surcoût administratif de l'index.

Requête	Sans Index (ms)	Avec Index (ms)	Gain (%)
Q1 Filmographie	1059.1	402.1	62.0%
Q2 Top Films	23.2	9.2	60.3%
Q3 Multi-rôles	219.1	235.3	-7.4%
Q4 Collaborations	344.5	0.3	99.9%
Q5 Genres Pop	35.3	43.3	-22.4%
Q6 Carrière	628.1	0.4	99.9%
Q7 Classement	80.0	89.0	-11.2%
Q8 Propulsés	3.4	3.9	-13.2%
Q9 Complexe	7.0	0.7	90.0%

TABLE 1 – Comparatif des temps d'exécution (Moyenne)

## 5 Conclusion Phase 1

Cette première phase a permis de valider la structure relationnelle des données. L'importation optimisée (transactions) et l'ajout d'index stratégiques ont rendu la base performante et prête pour la migration vers MongoDB (Phase 2).