



PROJET NETFLOOX

KAËLIG B.
LAURENCE B.
SAEEDULLAH R.Z.

Trello



Espace de travail Trello Gratuit

Buzzing Flash - Netfloo

Visible par l'espace de travail

Tableau

Planway Power-ups Automatisation Filtres

Partager

Tableaux

Membres

Paramètres d'espace de travail

Vues de l'espace de travail

Tableur

Calendrier

Vos tableaux

Buzzing Flash - Netfloo

Semaine 5 Fév

Semaine 12 Fév

Trello

Mise en place du projet

30 janv.

Créer les activités sous Trello.

30 janv.

Planning du projet

Connaissance des data

5 févr. - 13 févr.

Ajouter une carte

Github

Connecter les membres de l'équipe sur github.

6 févr. - 6 févr.

Ajouter une carte

Table SQL

Télécharger - Dézipper les fichiers

6 févr. - 6 févr.

Partage des fichiers

6 févr. - 6 févr.

Explorer les datas

Cette carte est un modèle.

2/3

Une base de données alimentées

Ajouter une carte

Interface web

Streamlit ?

16 févr.

Un modèle relationnel de données

Ajouter une carte

Un modèle relationnel de données

Explorer les datas

2/3

Créer base de données

<https://developer.imdb.com/non-commercial-datasets/>

Ajouter une carte

Python code

Un algorithme de prédiction de la popularité

Un algorithme de recommandation

Ouvrir sous python

31 janv.

Ajouter une carte



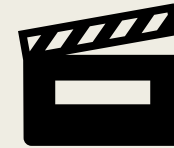
Ordre du jour :



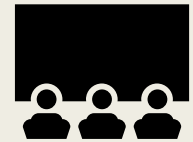
Introduction



Modèles



Résultats



Applications

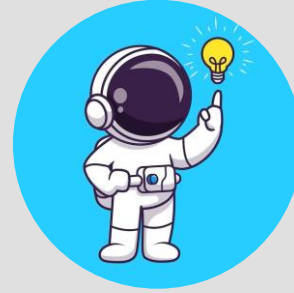
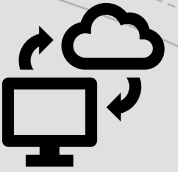


Les besoins :

Vous passez plus de temps
à choisir un film qu'à le
regarder ?



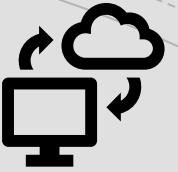
Les besoins :



Nous avons la solution !

Un service qui peut vous aider à faire le tri entre tous les contenus qui vous sont proposés et qui vous aidera à trouver les films et séries qui correspondent à vos critères.

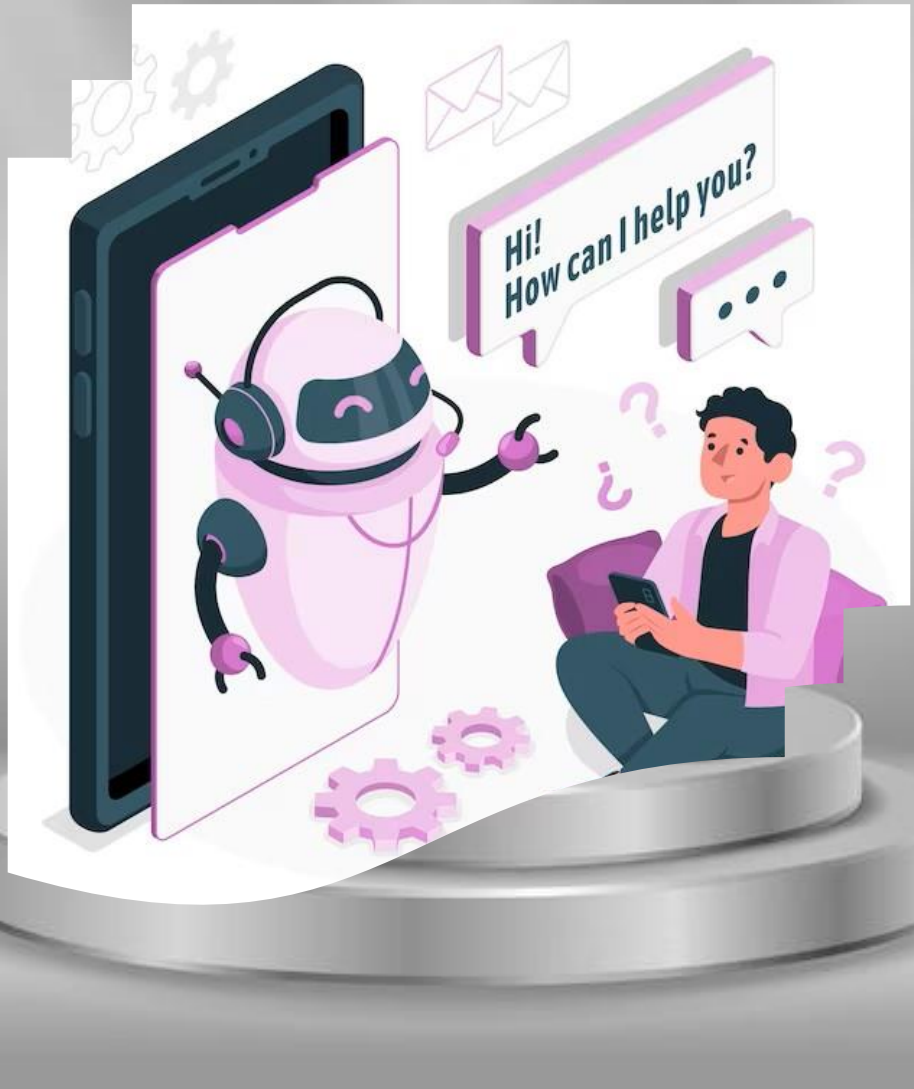
« Popcorn Movies »



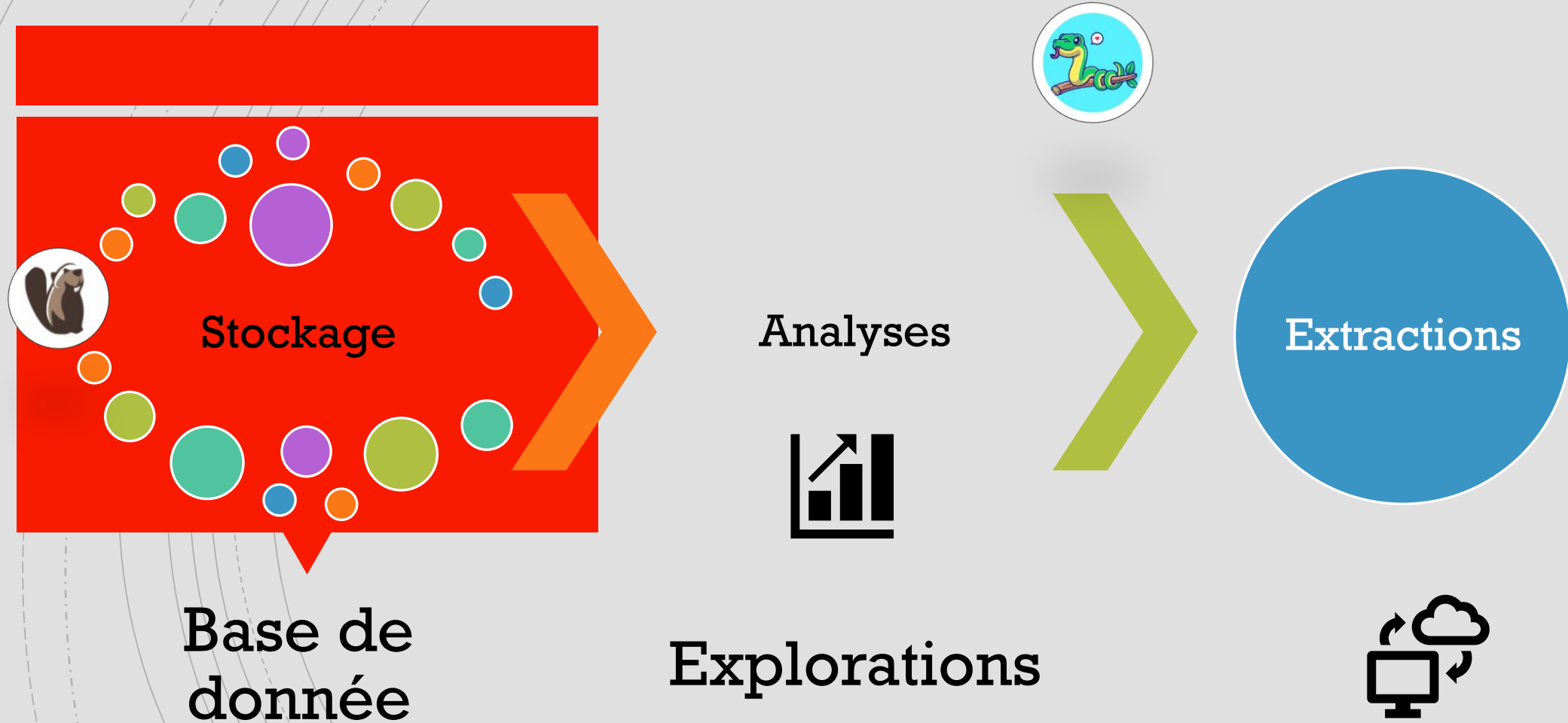
« Popcorn Movies »

Il s'agit tout simplement d'un service de **recommandation** de films qui prend en compte vos préférences en matière de films.

Recommandations personnalisées en fonction d'**un** film qui vous a plu.



Les données : 'Internet Movie Database (littéralement «Base de données cinématographiques d'Internet »).



Base de données :



It's
Time!



| public.title_episode | |
|----------------------|----------------------|
| tconst | varchar(50) NOT NULL |
| parenttconst | varchar(50) |
| seasonnumber | int4 |
| episodenumber | int4 |

| public.title_ratings | |
|----------------------|----------------------|
| tconst | varchar(10) NOT NULL |
| averagerating | float4 |
| numvotes | int4 |

| public.title_akas | |
|-------------------|----------------------|
| tconst | varchar(11) NOT NULL |
| ordering | int4 NOT NULL |
| title | varchar(600) |
| region | varchar(5) |
| language | varchar(5) |
| types | varchar(20) |
| attributes | varchar(60) |
| isoriginaltitle | int4 |

| public.recommendations | |
|------------------------|------------------|
| id | serial4 NOT NULL |
| movie_name | varchar(255) |
| recommended_movie | varchar(255) |

| public.title_principals | |
|-------------------------|----------------------|
| tconst | varchar(11) NOT NULL |
| ordering | int4 NOT NULL |
| nconst | varchar(11) |
| category | varchar(25) |
| job | varchar(40) |
| characters | varchar(85) |

| public.title_basics | |
|---------------------|--------------|
| tconst | varchar(11) |
| titletype | varchar(500) |
| primarytitle | varchar(500) |
| originaltitle | varchar(250) |
| isadult | int2 |
| startyear | int2 |
| endyear | int2 |
| runtimeinminutes | int4 |
| genres | text |
| genres_type | _varchar |

| public.name_basics | |
|--------------------|----------------------|
| nconst | varchar(11) NOT NULL |
| primaryname | varchar(50) |
| birthyear | int2 |
| deathyear | int2 |
| primaryprofession | varchar(50) |
| knownfortitles | text |
| profession | _text |
| connu_pour | _text |

| public.title_crew | |
|-------------------|----------------------|
| tconst | varchar(20) NOT NULL |
| directors | text |
| writers | text |

| public.view_fin | |
|------------------|--------------|
| directors | _varchar |
| writers | _varchar |
| genres_type | _varchar |
| titletype | varchar(500) |
| runtimeinminutes | int4 |
| tconst | varchar(11) |
| language | _text |
| actor | _text |
| originaltitle | varchar(250) |
| averagerating | float4 |

| public.nbr_language | |
|---------------------|------------|
| language | varchar(5) |
| nbr | int8 |

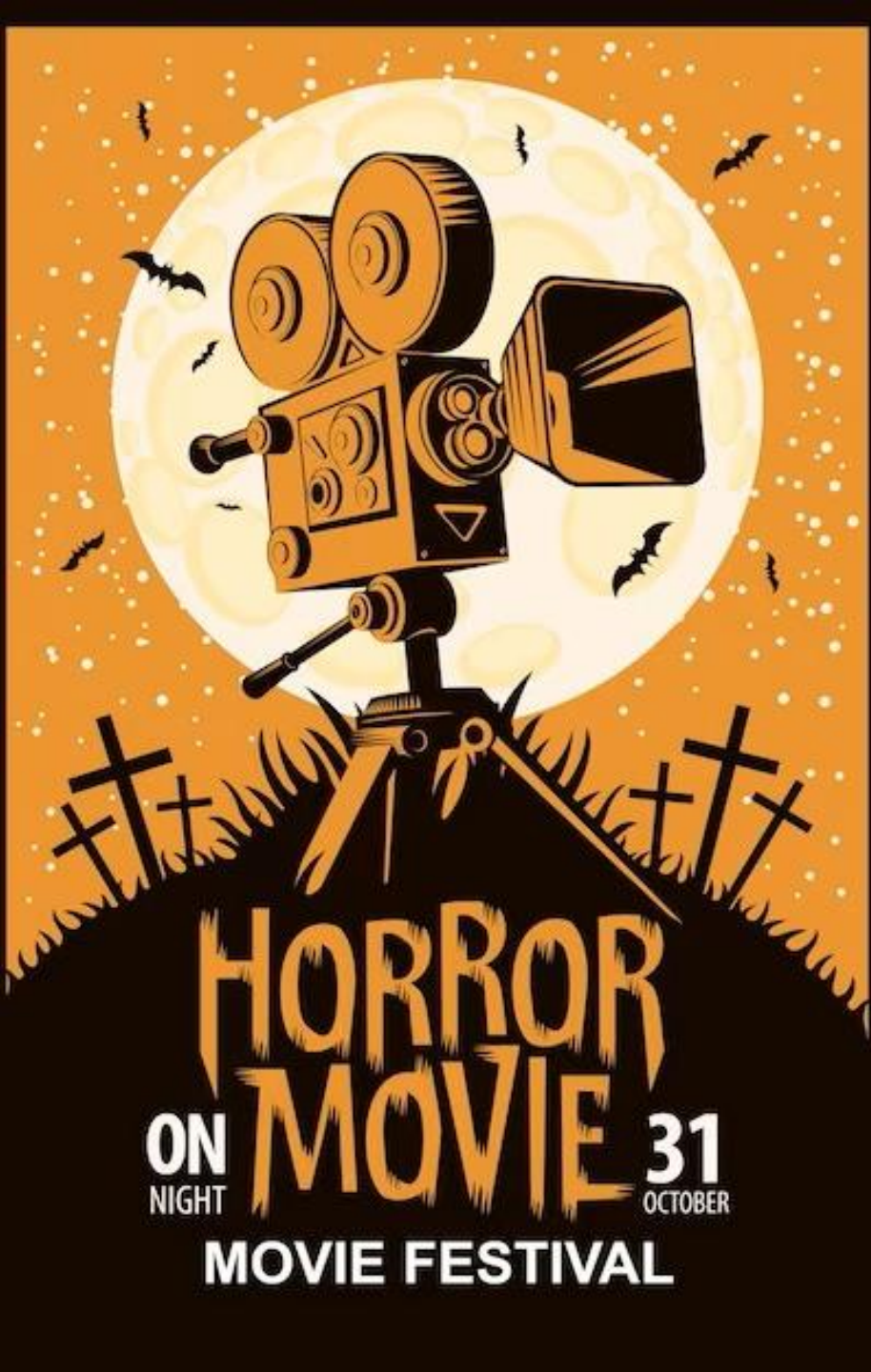
| public.nbr_region | |
|-------------------|------------|
| region | varchar(5) |
| nbr | int8 |

| public.title_crew_view | |
|------------------------|-------------|
| tconst | varchar(20) |
| directors | _varchar |
| writers | _varchar |

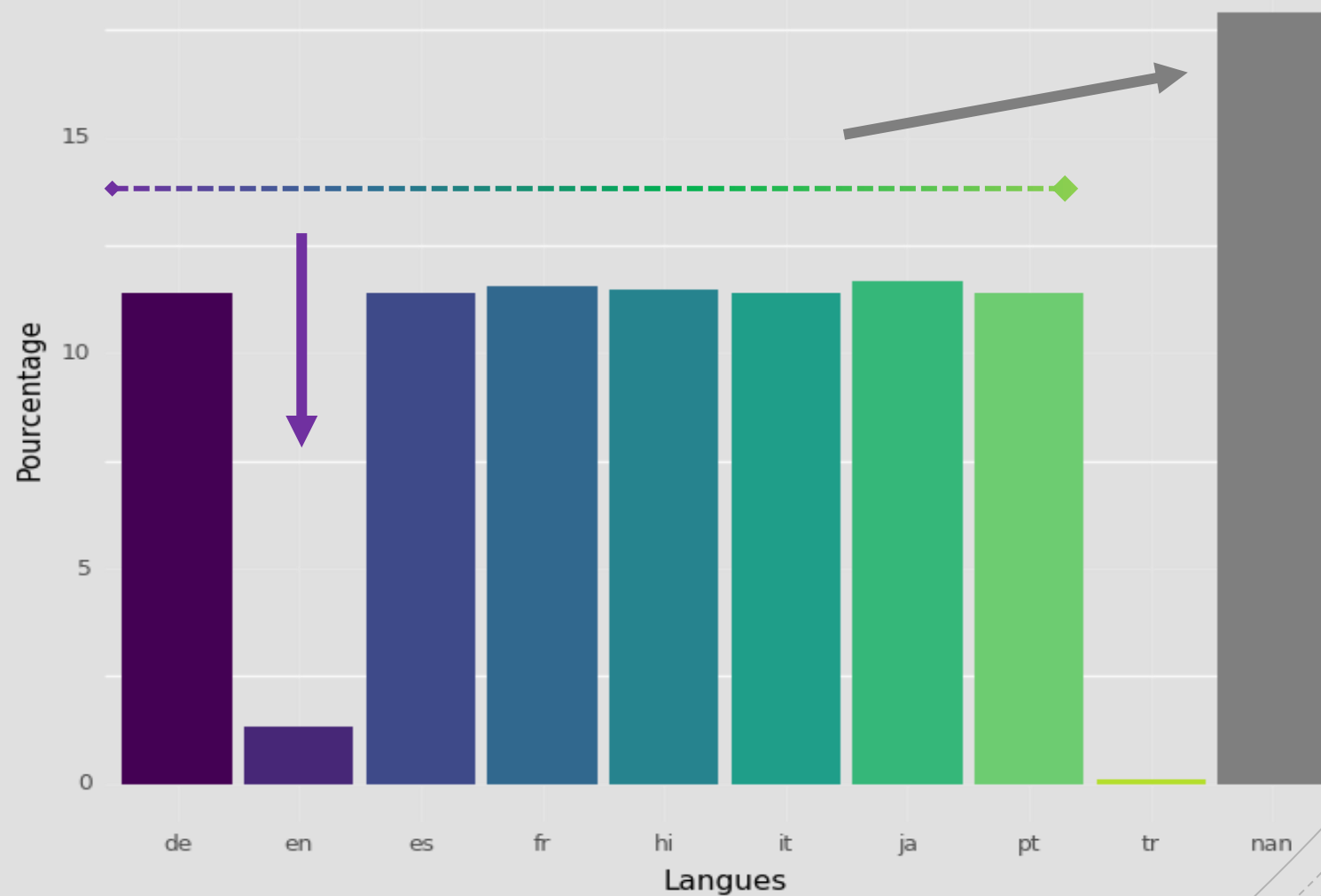
| public.vote_etoile_films | |
|--------------------------|--------------|
| numvotes | int4 |
| averagerating | float4 |
| primarytitle | varchar(500) |

| public.acteur | |
|---------------|-------------|
| actor | _text |
| tconst | varchar(11) |

| public.category_job | |
|---------------------|-------------|
| category | varchar(25) |
| compte | int8 |



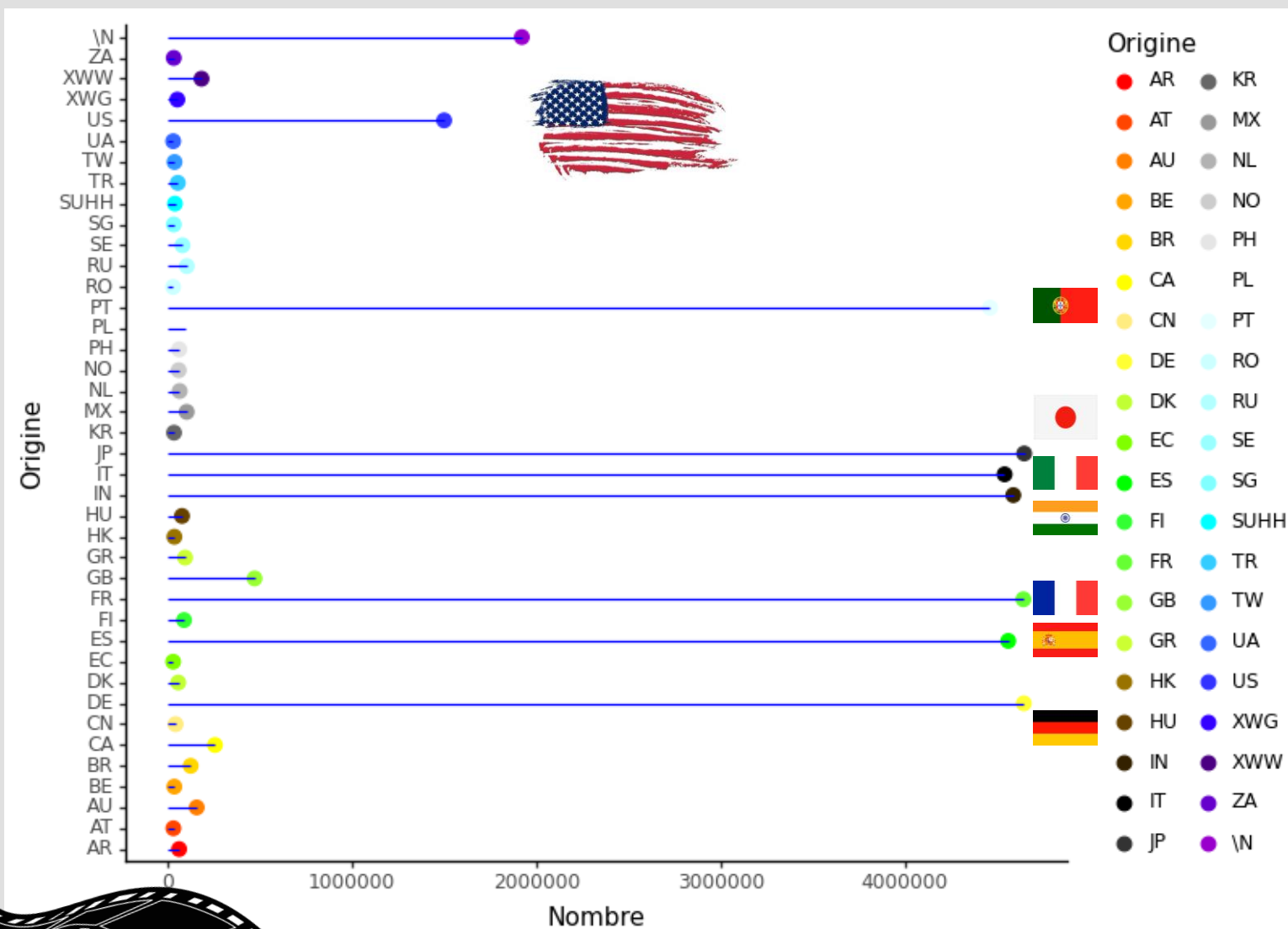
Top 10 des langues pour l'ensemble de la BD



107 langues

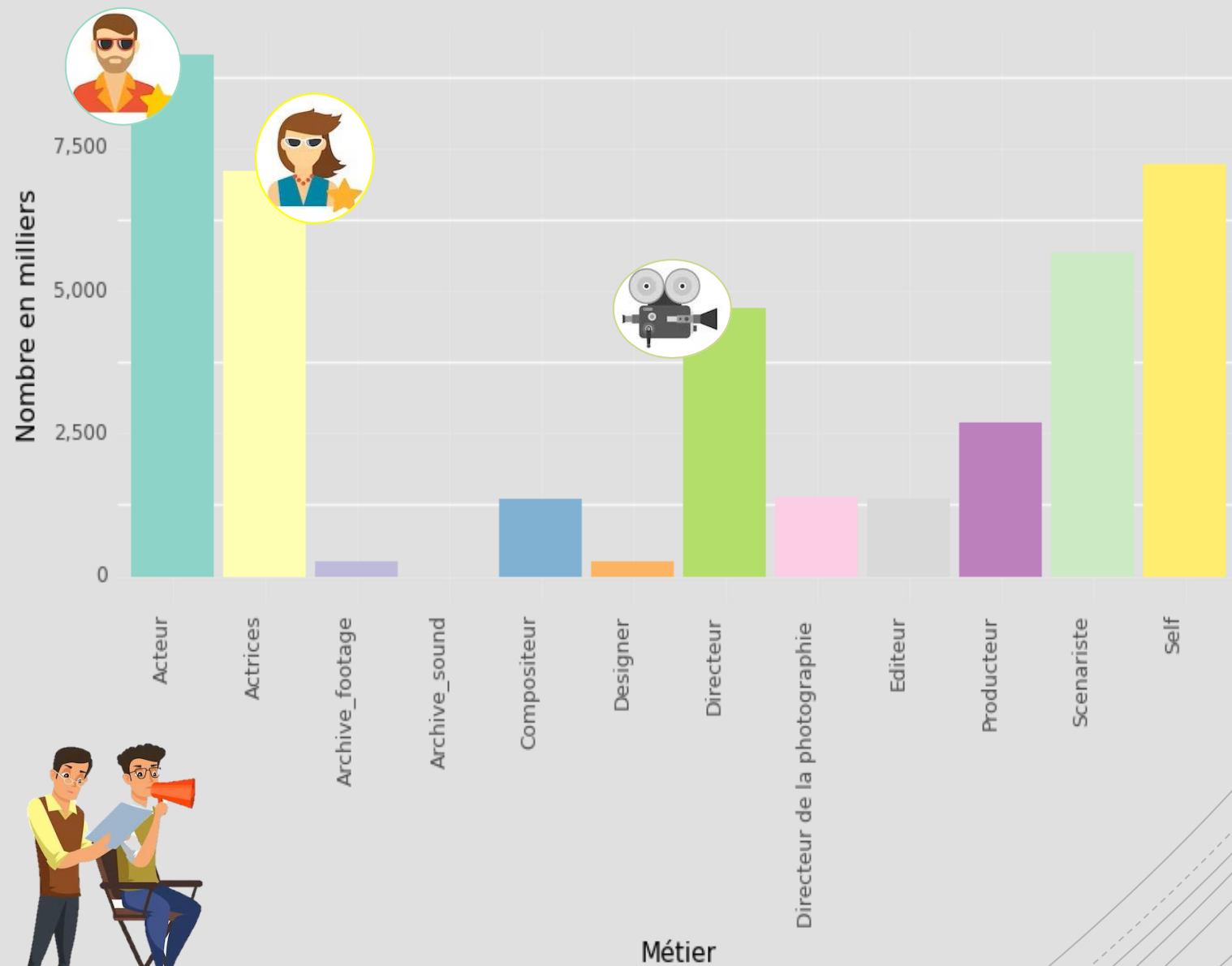


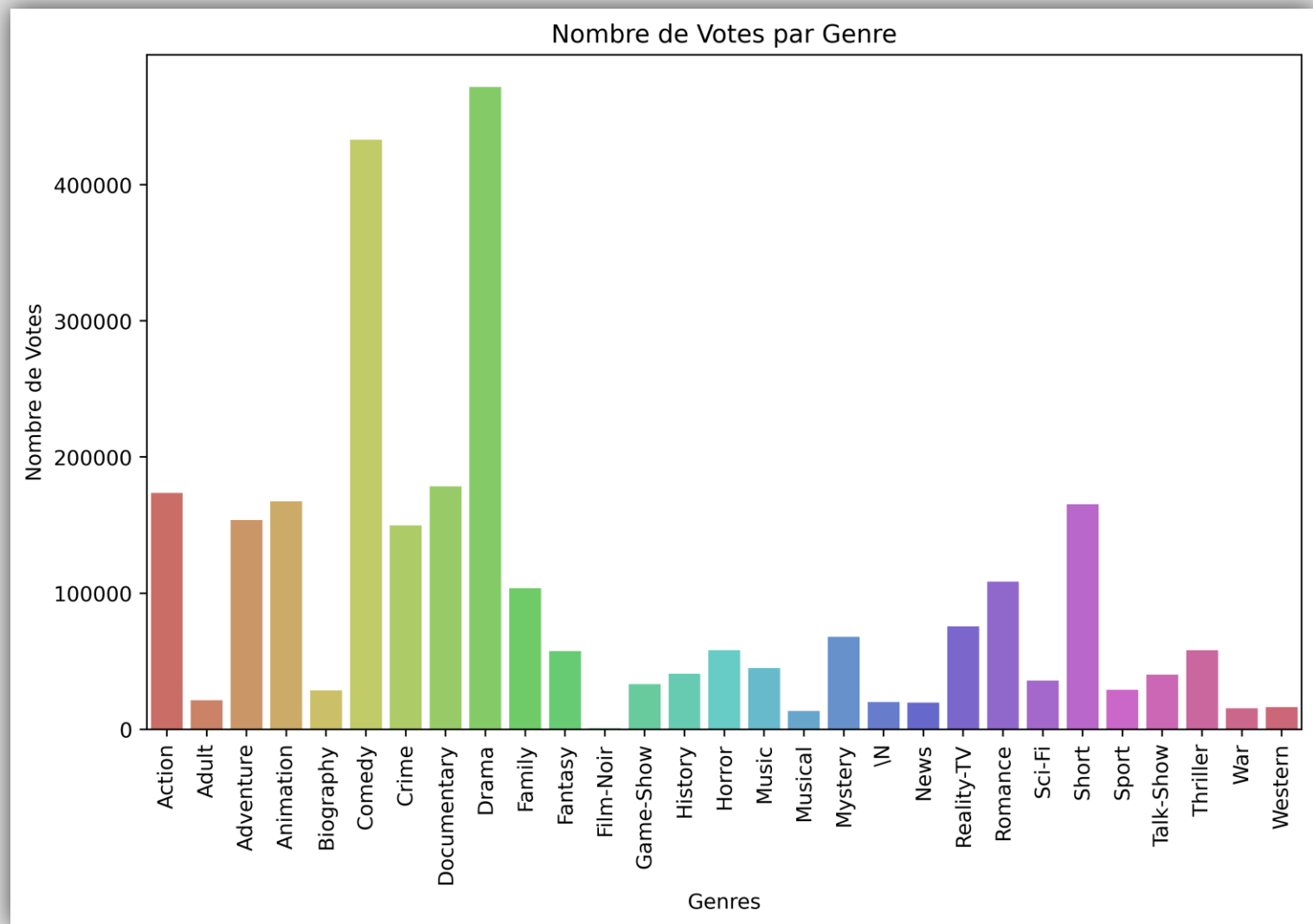
Top 40 de l'origine pour l'ensemble de la BD





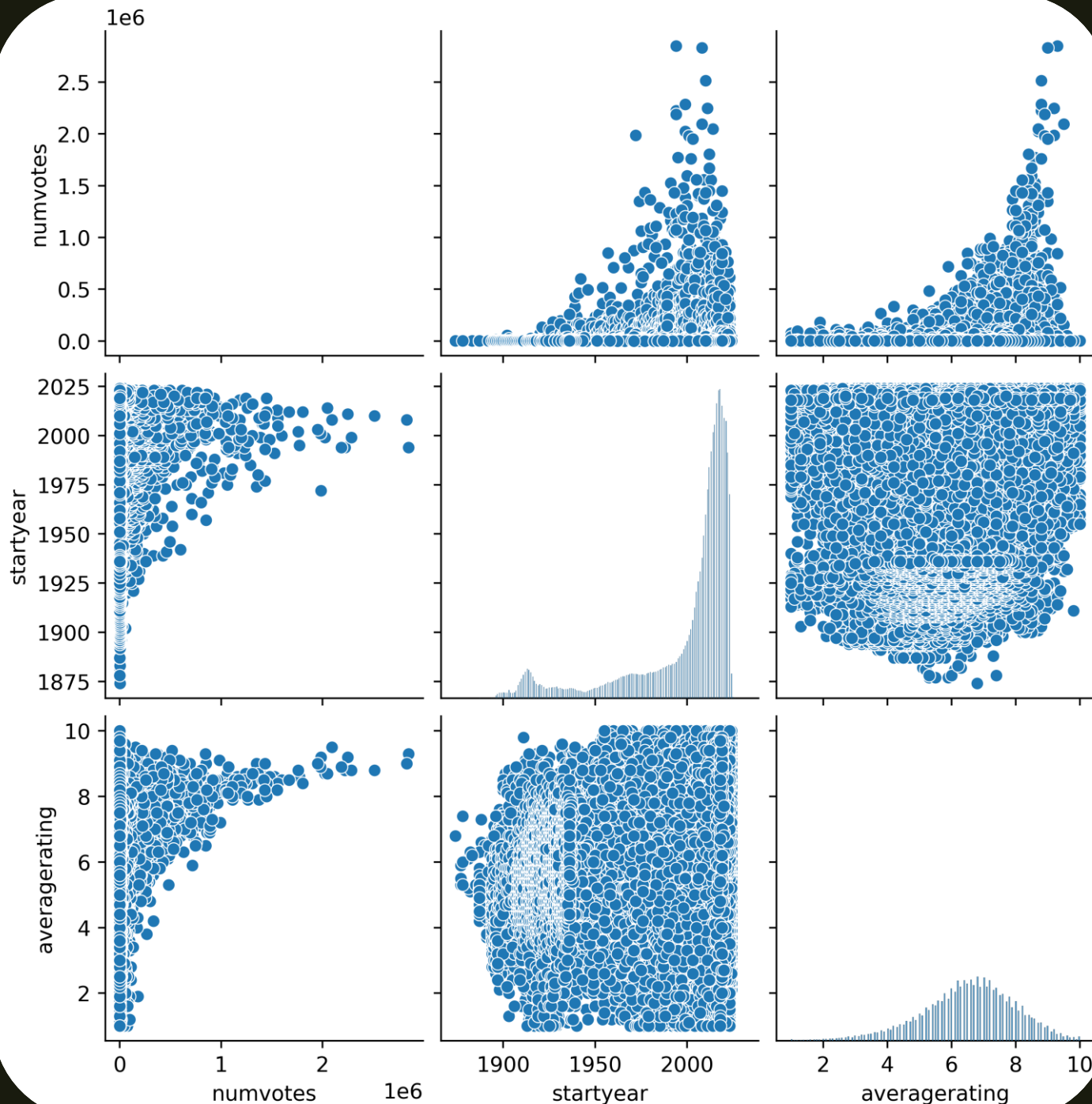
Nombre de personnes par catégorie de métier







Relation entre la notation, le nombre de vote et les années de production.



Modèle :



BDD

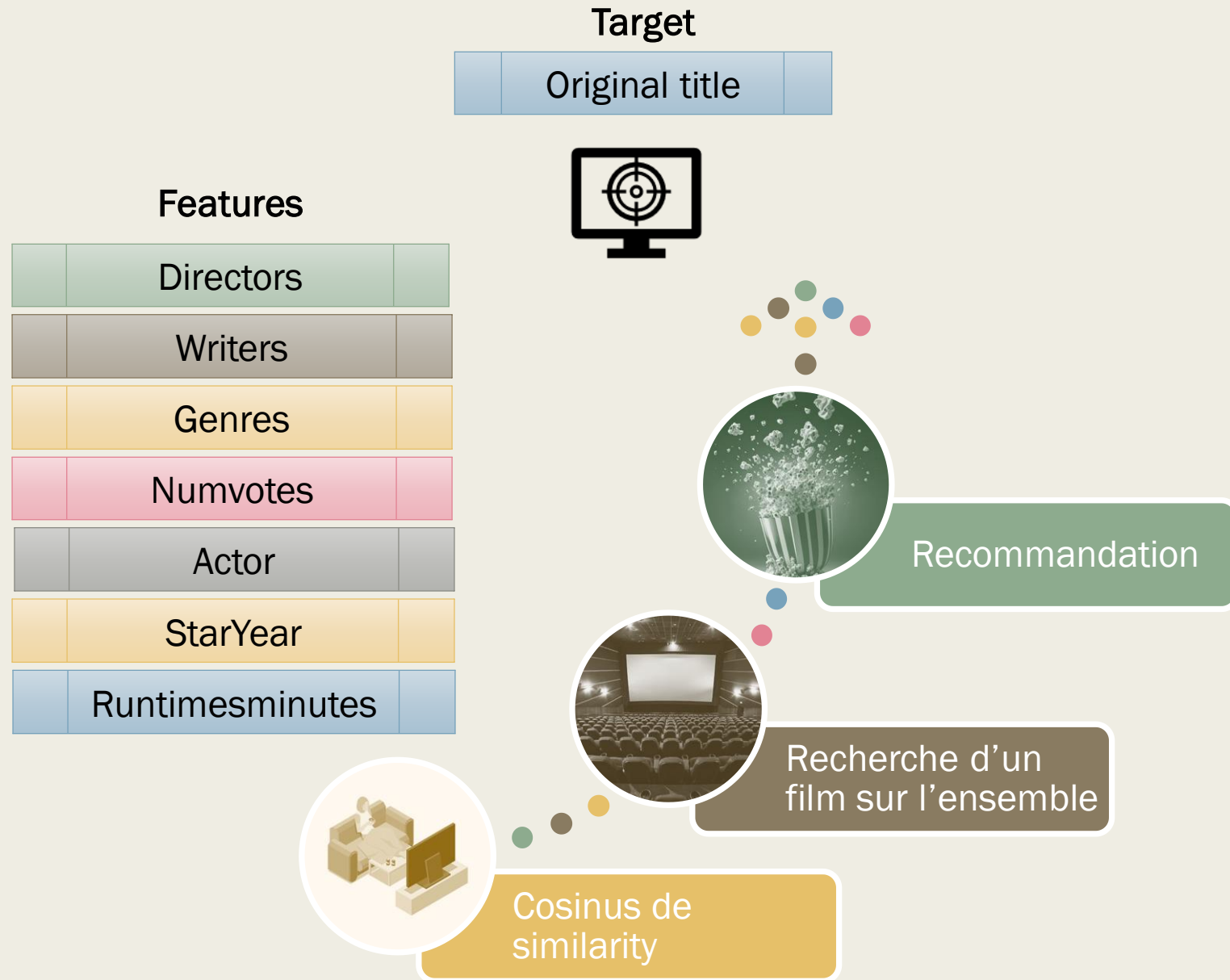
Extraction CSV

Analyses

Cosinus de similarité



Modèle :



Modèle de recommandation :



- Possible sur nos ordinateurs



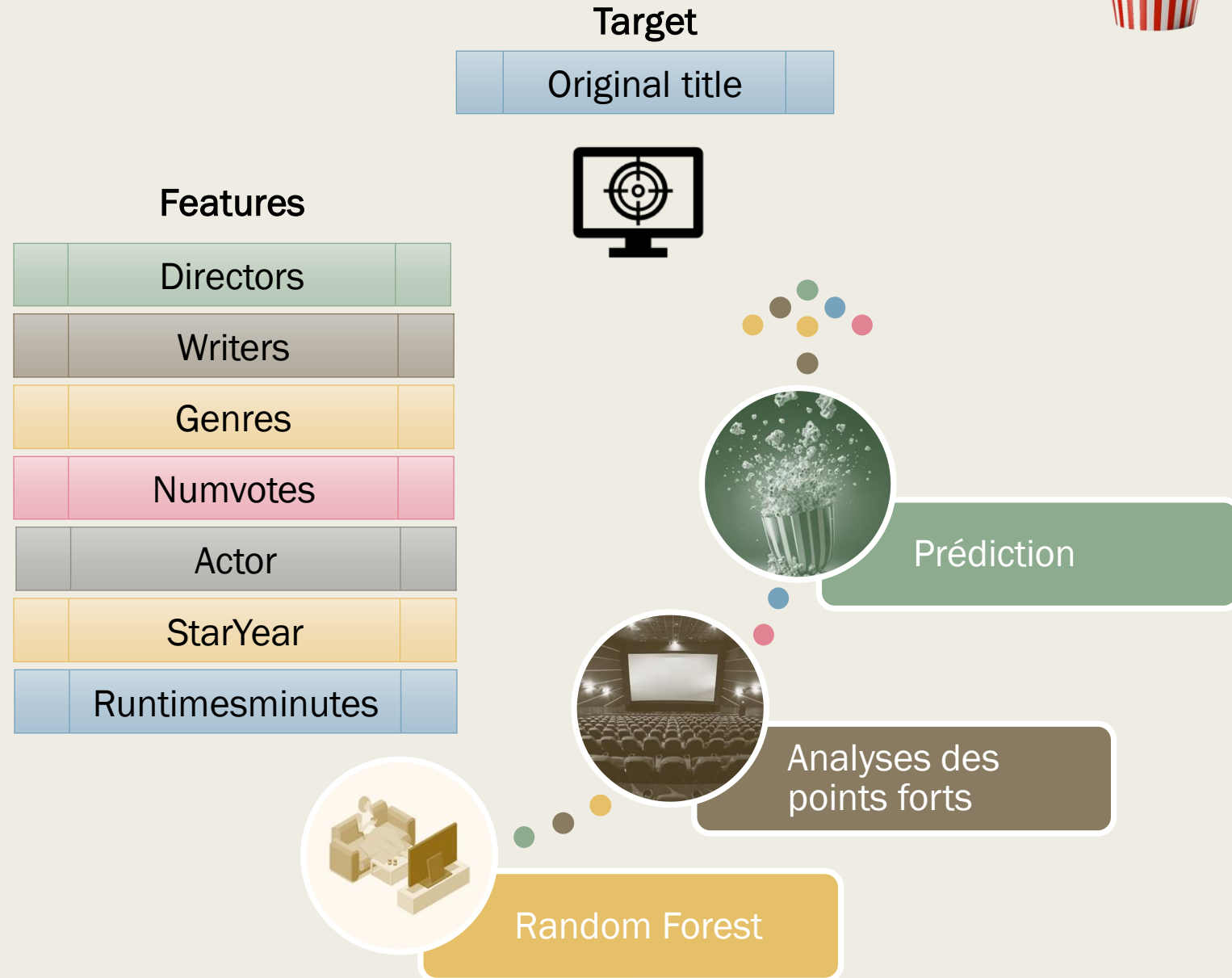
- Re-calculer le cosinus à chaque film



- Pour éviter une matrice trop importante
- 5 films les + similaires

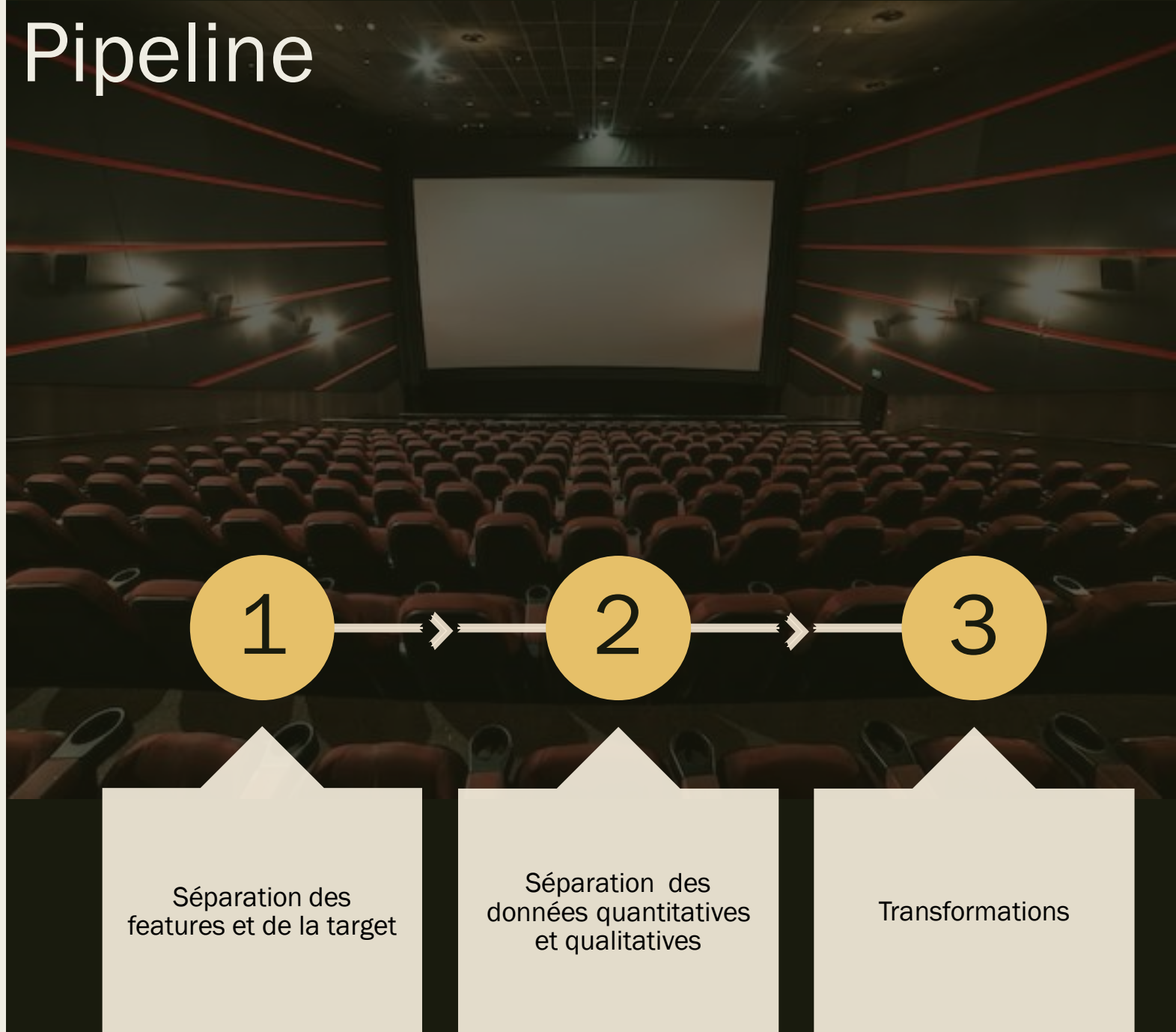


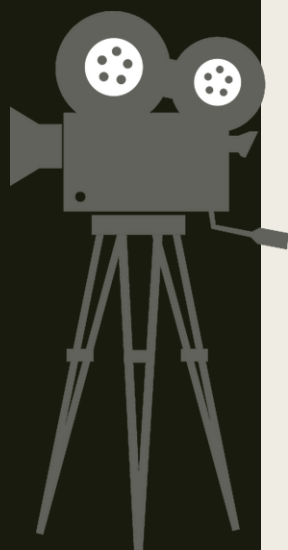
Modèle :



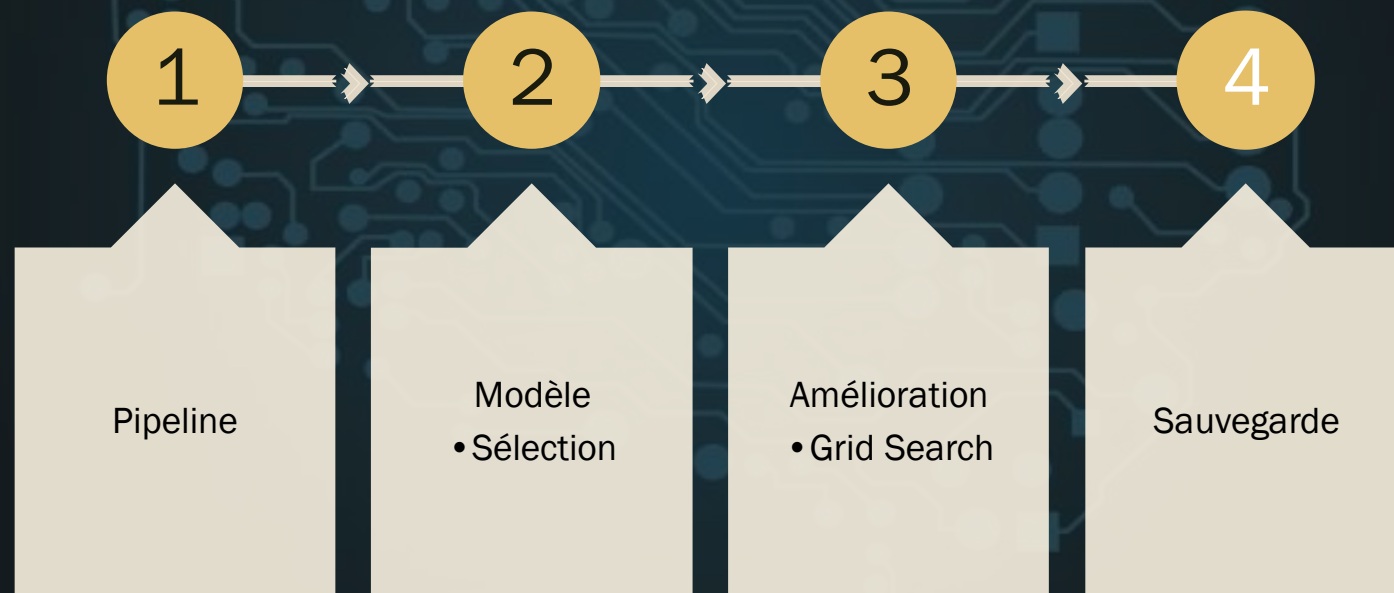


Pipeline





Comment entrainer le modèle



Modèle :



```
def pipeline(model=DecisionTreeRegressor(), num=[SimpleImputer(),StandardScaler()],
            text=[OneHotEncoder(handle_unknown='ignore')]):
    colonne_numerique=['startyear','runtime','numvotes']
    colonne_nominal=[colonne for colonne in X.columns if colonne not in colonne_numerique]

    pip_num = Pipeline(steps=[(f"etape_num_{count}", i) for count, i in enumerate(num)])
    pip_text = Pipeline(steps=[(f"etape_text_{count}", i) for count, i in enumerate(text)])

    preprocess = ColumnTransformer(transformers=[('text', pip_text, colonne_nominal),
                                                ('num', pip_num, colonne_numerique)])

    model = Pipeline(steps=[('preprocessor', preprocess), ('model', model)])
    return model
```

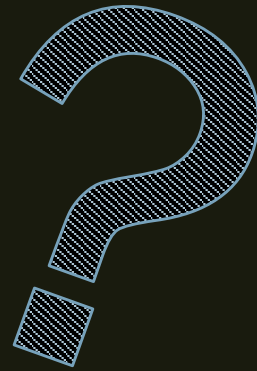
```
def hgs(df,param,cv):
    y = df['averagerating']
    X = df.drop(['averagerating','titletype','features'],axis=1) #df['message']

    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
    grid = HalvingGridSearchCV(estimator=pipeline(),param_grid=param,scoring='neg_mean_squared_error',
                              cv=cv,n_jobs=-1, verbose=1, error_score="raise")

    grid.fit(X_train, y_train)

    best_score = grid.best_score_
    best_params = grid.best_params_
    training_time = grid.cv_results_['mean_fit_time'].mean()
    return({'best_score': best_score,
            'best_params': best_params,
            'training_time': training_time,
            'fitted_model': grid.best_estimator_})
```


On regarde
quoi ce soir



Application



© 2006 The Authors
Journal compilation © 2006 Blackwell Publishing Ltd

