



# CaeaFlix project

*Algorithme de recommandation de films*

# Sommaire

- **Qui somme-nous ?**
- **CaeaFlix ... c'est quoi ?**
  - L'objectif du projet
  - Quels films ?
  - Constitution de la BDD
  - Data viz
  - Modèle
  - Modèle final
- **Démo**
- **Limites**
- **Recommandations**
- **Next Steps**

# Qui sommes-nous ?

# Equipe

## Scrum master :

Arnaud Le Naourès

## Product Owner :

César Ozeer

## Team Members :

Antoine Rafflegeau

Esrin Erdem

# CaeaFlix ... C'est quoi ?

# L'objectif du projet

## Projet Caeaflix :

- Relancer les ventes d'un cinéma dans la Creuse

## Demande :

- Mettre en place un système de recommandations de films

# Quels films ?

## Approche business :

- S'appuyer sur une étude consommateurs afin de prendre les bonnes décisions

## Sélection films:

- Films(movies)
- diffusé en français ou en anglais (*mais seulement pour ceux aux USA*) ...
- ... dans plus de 7 pays
- note moyenne > Q1
- nombre de votes > Q1

La base de données finale contient près de **15 000 films**.

# Constitution de la base de données

## Données de départ :

- *title\_akas, title\_basics, title\_principals, title\_ratings, name\_basics*  
⇒ *Fusion des tables*

## Manque d'informations :

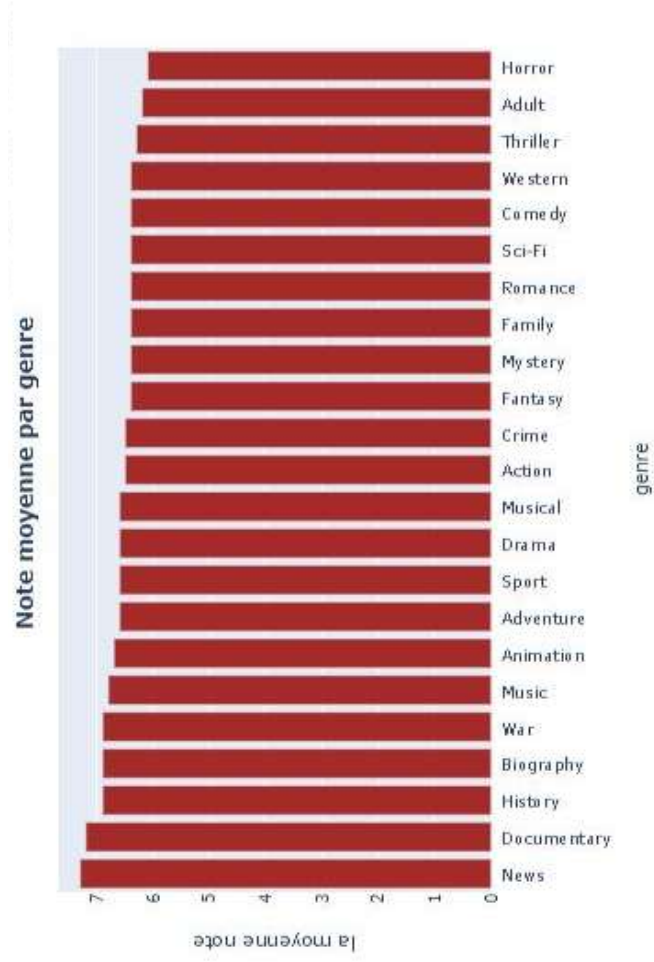
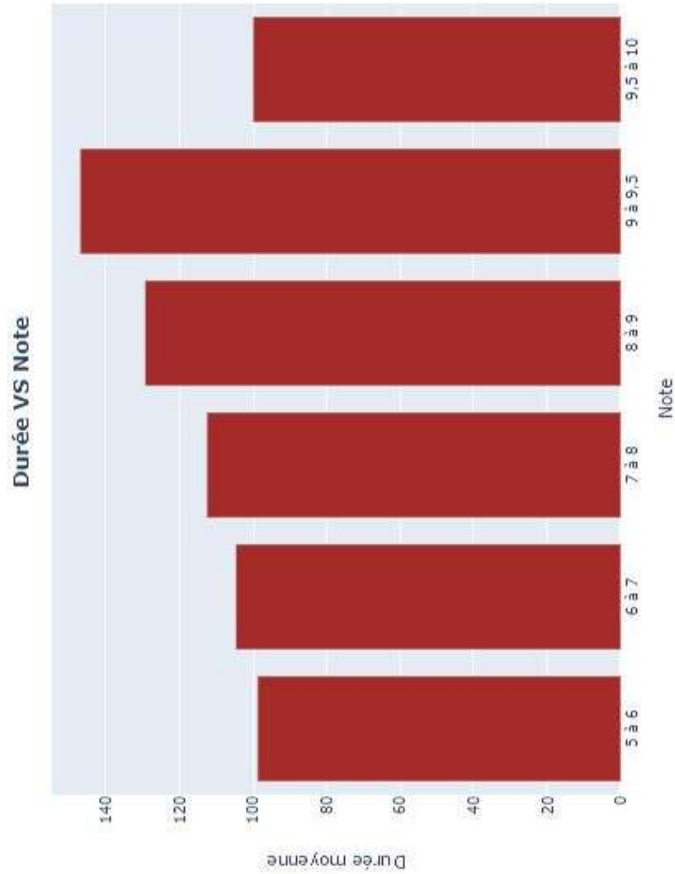
- recettes, budget, affiche du film, mots-clés, lien vers le site du film  
⇒ recours à des API pour récolter ces informations

## Data cleaning :

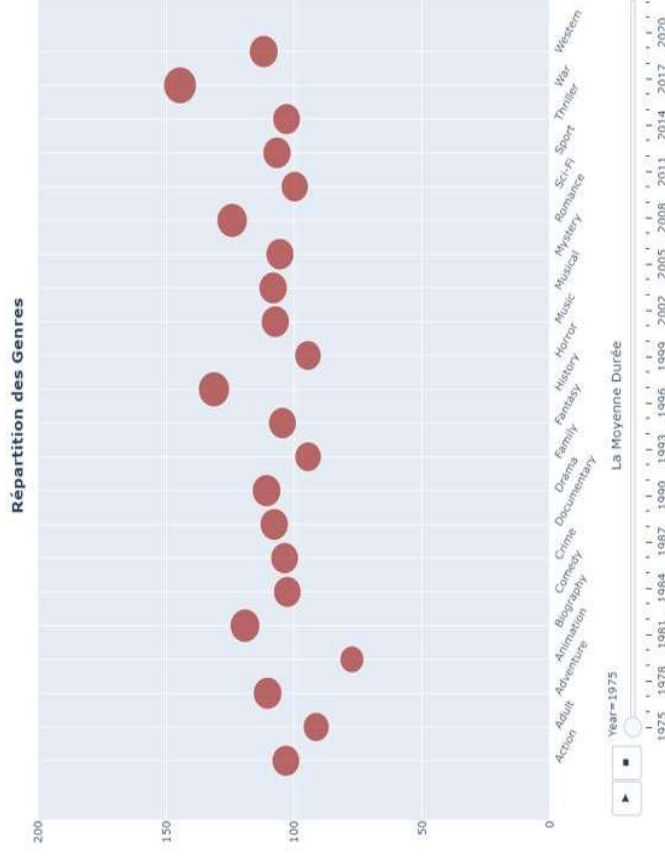
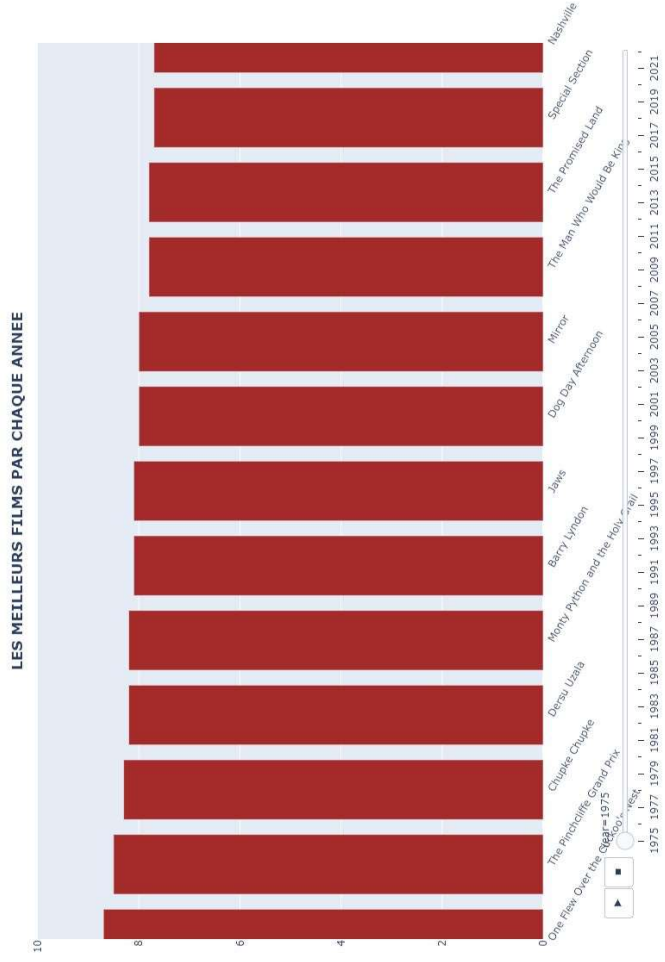
- Regroupement avec une ligne par film et une colonne par catégorie (acteur, actrice, réalisateur)
- Suppression des "NaN"



# Data viz (1/2)



# Data viz (2/2)



# Modèle

## Modèle : K-Nearest Neighbors (K-NN)

- Variables  $\Rightarrow$  positionnent les films les uns par rapport aux autres
- Recommandation  $\Rightarrow$  film le plus proche en tenant compte de toutes les variables

# Modèle final <sup>(1/2)</sup>

Mise en place de 4 dataframes :

- Films “anciens” (avant 1999 compris)
- Films “récents” (après 1999)
- Anime (seulement les anime)
- Indiens (seulement les films de Bollywood)

Création de variables :

<u>Binaires</u>	<u>Scaled</u>
Genre	Runtime minutes
Key words	Average rating
Pays / ISO / nom de la (des) société(s) de production	Number of votes
Nom réalisateur / acteurs / actrices	Start year
	Diffusion

# Modèle final <sup>(2/2)</sup>

- **3 modèles :**

- Modèle général (films les plus similaires) :
  - X = Keywords + Nom des producteurs + Nom des pays des producteurs + Nom des ISO des pays des producteurs + Nom réalisateur + Nb diffusion + Genres + Langue
- Modèle par réalisateur :
  - X = Nom réalisateur + note moyenne
- Modèle par studio de production :
  - X = Nom producteurs

# Démo

# Limites

# Limites

## Base de données trop lourde (30K+ colonnes) :

⇒ Lenteurs dans l'exécution du code

## Volonté de donner un poids différents aux variables :

⇒ Le modèle KNN tel qu'on l'utilise ne le permet pas



# Recommendations

# Recommandations

- **Trouver d'autres modèles** afin de ne pas avoir autant de colonnes et pouvoir donner plus d'importance à certaines variables  
⇒ un modèle possible serait le Scipy scalar (hot encoding)
- **Utiliser la méthode GridSearch / Randomized Search**

# Next Steps

# Next Steps

## Pour apporter plus de services:

- Création d'un système de recommandation fondé sur la consommation des autres clients ⇒ User based
- Création d'un système qui récupère les recommandations faites aux utilisateurs et donne des indications sur les films à prochainement mettre à l'affiche
- Création d'un système de recommandation fondé sur les deux systèmes précédents ⇒ Hybrid based
- SMS des recommandations de films après avoir acheté une place pour une séance au cinéma.



# CaeaFlix project

*Algorithme de recommandation de films*