

Topic Modelling of Movie User Reviews

- SUBMITTED BY :
 - TANVI ARORA

User Reviews for Movies

- Source : imdb.com
- Genre : Biography
 - The Irishman
 - Bohemian Rhapsody
 - When They See Us
 - RocketMan
- What data is available on imdb ?
 - User Rating (scale of 0 to 10) - structured data
 - (eg : positive reviews : >5 , negative reviews <5)
 - User reviews (ranging from 0 to 1000s) - unstructured data
 - Votes

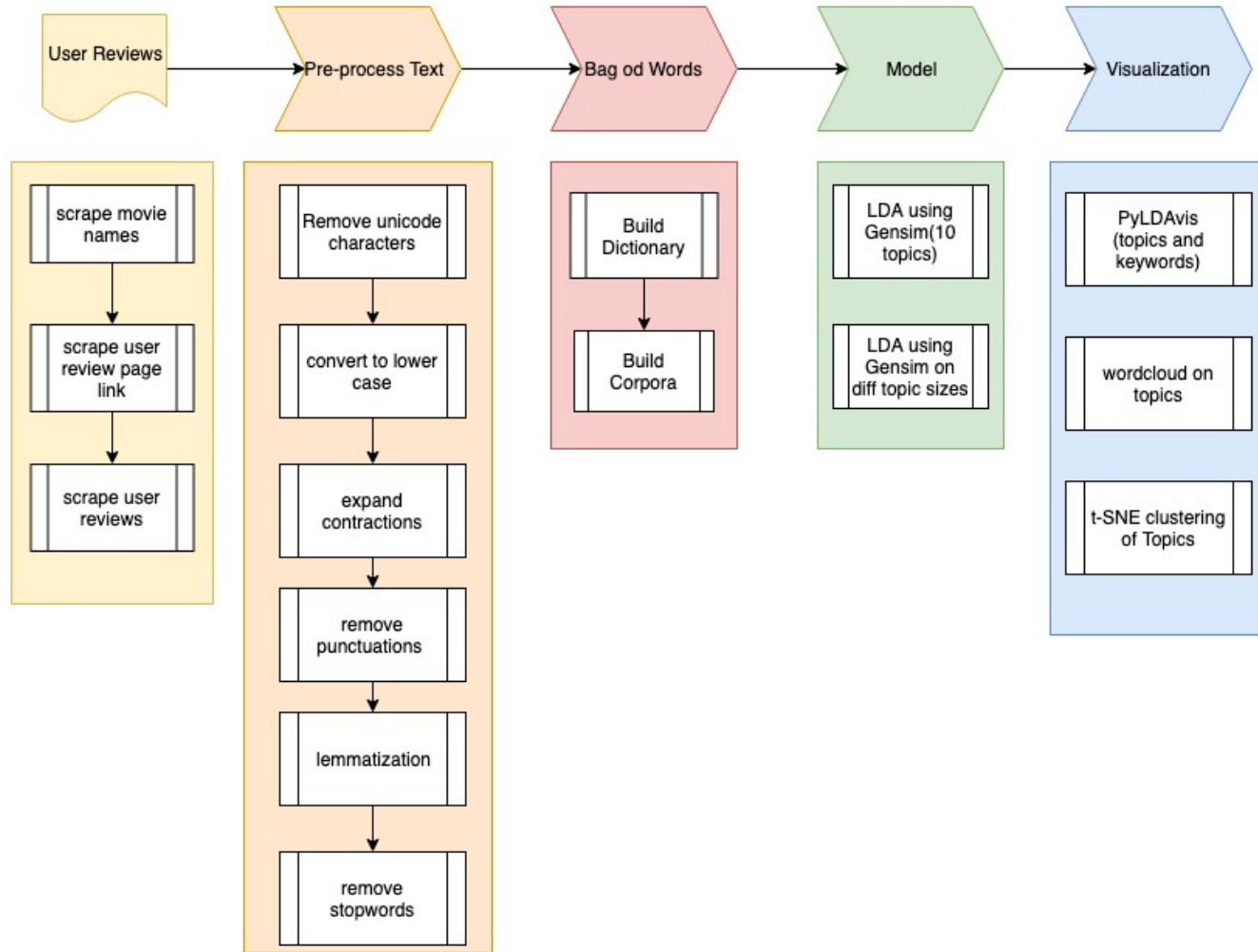


What are they talking about ??

- User reviews can be positive or negative based on the user rating
- But really what are the users talking about ?
 - Movie
 - Actor
 - Performance
 - Type of story
- How do I find that ??
 - Topic Modeling
 - Is a type of statistical modeling for discovering the abstract topics that occur in a collection of documents
 - Help in making sense of the unstructured data

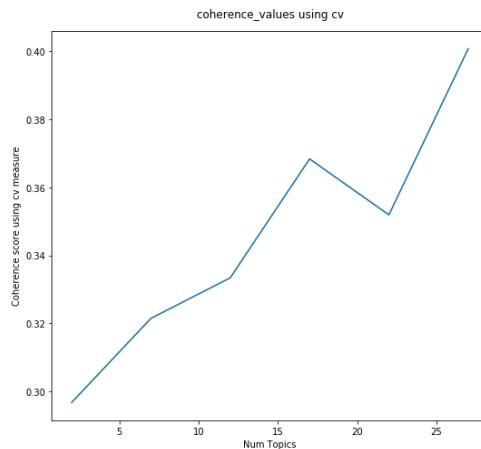
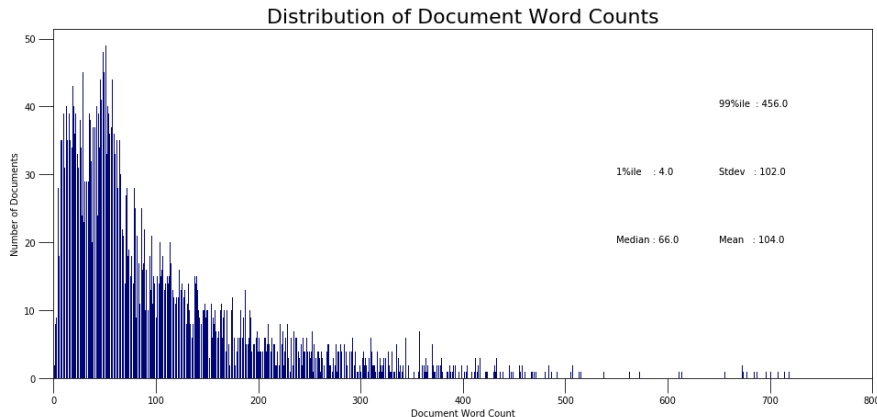


Pipeline

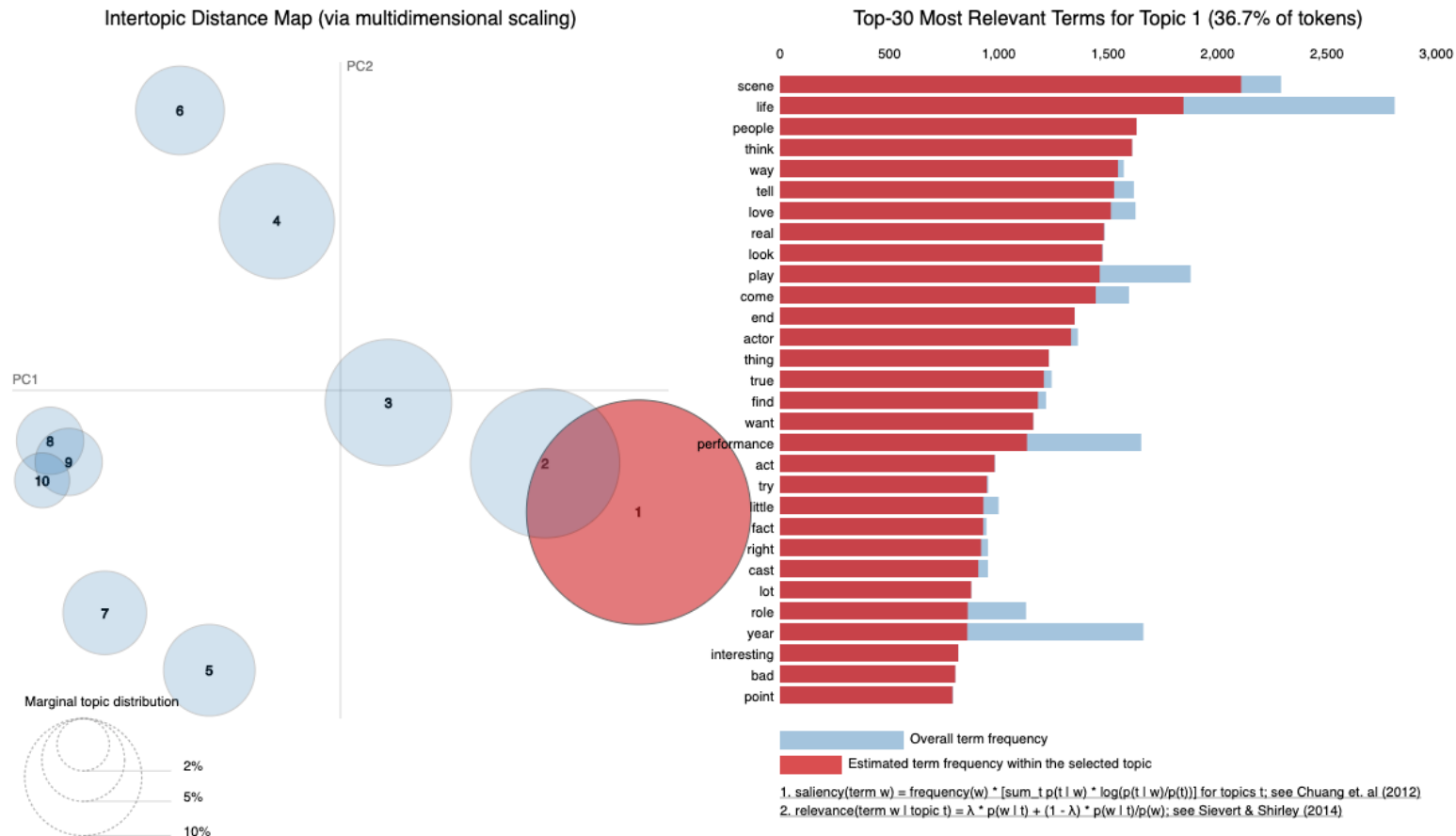


Model

- How much data are we talking about ~ 4500 user reviews
- Latent Dirichlet Allocation (LDA)
- A few open libraries exist but main contender is Gensim
- Cannot automatically determine number of topics
- Tested with a range of 2 to 30 with an interval of 5 for the number of topics
- Topic Models give no guarantee on the interpretability of the output
- Coherence Score : measures the relative distance between words within a topic
 - We want coherence score to be as high as possible



Topic and Keywords



T
O
P
I
C
S

Topic 0
piece jackie
old paul knight
william
accent brian
christopher
version

Topic 2
war life
year dream
work
script action
performance
man
world

Black
man &
woman

Topic 4
lizzie
ip
black
woman
yen
murder
art
tennis
borden

Shakespear
-en era

Historic
al wars
&
politics

About
the
movie

Topic 6
people scene
way real
love look
think life
play
tell

Topic 8
relationship
rise family
wilson
period life son
boy
child wind

Topic 1
branagh
selma tyson
escobar webb
marilyn dench pablo
kill
civil

Topic 3
shakespeare
terrible society jump drug
king
abdul queen
philippe victoria

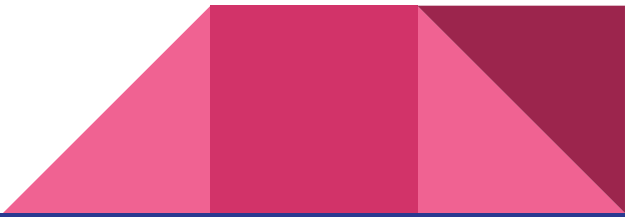
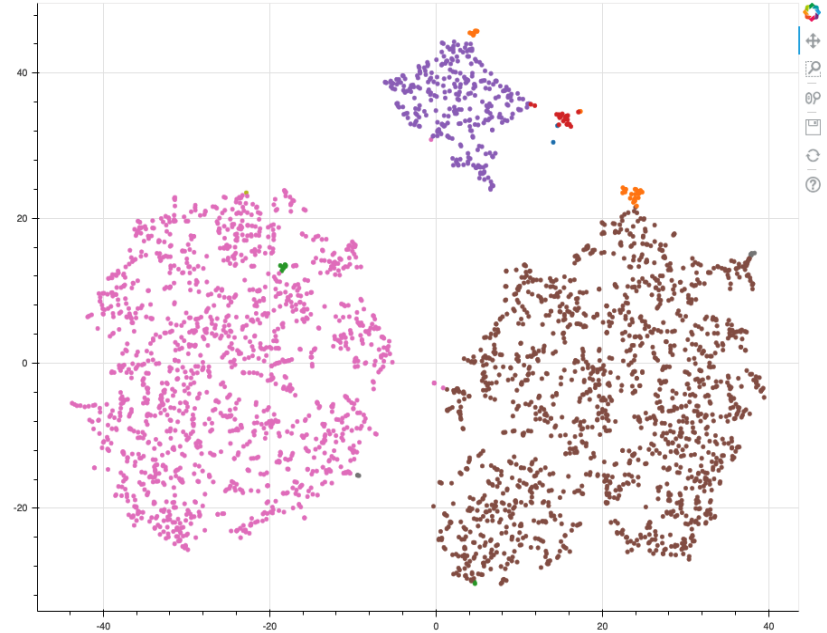
Topic 5
historical
political
group
mary uses
history
wire american
truth war

Topic 7
actress
greene ray
linda
sevigny michael
prison renner
jeremy lovelace

Topic 9
danny
ground hang self
pacino
irish
miyazaki
fail jiro
decent

International
movies

t-SNE Clustering of 10 LDA Topics





Thank you