NLP Task: part of speech Tagging and information situitival

Alh!

By using NLP libraries to perform pours

PROURAM CODE :

import spacy

nep = spany. road (" en- con- web, sm")

teact = "A1- ouriver platforms personalize rearning

paths and help student grasp

conteple"

doc = rup (tepit)

for token in/doc:

print (] "{toren lest: 15} >
{toren pos_ 3"}

Jeon /sicleaun jeatur eschartion tour import

Jon Wiceaun militis. pourulse import losing.

dovumen 13 = [

provide ear time feedback!

a findelligent fotosing system odapt to each student's reduning style's

2025/10/25 1

Al -> propri PUNCT driven -> VERB platjosms -> NOUN personalize -> VERB learning -> VERB paths -> NOUN and -> CONJ help -> VERB Students -> NOUN grasp -> VEKB Contepls -> NOUN fastes -> Apu

-> FUNCT

query = "How does At support students in learning?"

conpus = documents + cavery

vectonizer = That vectonizer()

ftidt_mainir = vectonizer. filtransjoum (conpus)

Similarities = 108 in- similarity (cticy - matrix [-1],

Cticy - matrix [:-1] · flatten()

nanted dous = sontea (zipl similarities,

documents 1, kivese = TRUE)

print (" n Top revant documents: \n")

Jon slore, doc in nanted_docs:

print (] " store : 4 store : . 33 -> {doc 30.

Thus rup peyermanio has been essewled successfully.

2025/10/

Exno:6 Exploratory Data Analysis with Python

AIH!

To do exploratory data Analysis with python

Program:

emposit pandas as pol
emposit numpy as np
emposit matplotlib pyplot as plt
emposit seabarn as sos

dy = pd. read_csv (4 retflix_ titles.csv")

print (dy injou)

print (dy. head ())

print (dy discribe (melled , 'all'))

print ("Number of unique countries", of ['country'].

posint ("Number of unique directors", dy ['director].

nunique)

print (dy ['release_year']. value_counts())

print (dy ['release_year']. value_counts(). head())

print(dy group by (['(ountry', '(type)]). size().

Bost values Castending = Palee).

nead (10))

of ['dak_added'] = pd. to_datetime (of C'dak_add

gormat = 'misted', exors = 'Co

2 class 'pandas · core · frame · Data Frame'>
Range Index = 8807 entries, 0 to 8806

Data columns (fotal 12 columns)!

column Non-Null Count Dtype 0 show_id 8807 non-null object 1 type 8007 non-null object 2 title 8807 non-null object 3 director 6173 non-null object 4 cast 7002 non-null object country 7976 non-null Object date-added 8797 non-null object 7 Micase_year 3807 nonnull int 64 nating 8803 non-null object

drypes: Int 64 (1), object (u)

Number of unique countries: 748

Number of unique directors: 4528

Type

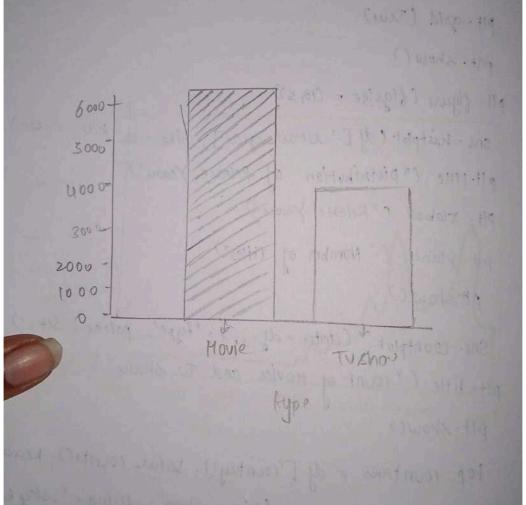
Hovie 6121

TV show 2676

Moder : count, drype 1 Int 64

```
dy set induc ('date_added', inplace = Tim)
monthly - content = dy nesample (41) - sige ()
 plt-ligue (figsige = (12,6))
  monthly - content · plot ()
  plt-title ("Netflish Content Added over Time")
  plf. x label ("pale")
    plf- glabel ("Number of Titles Added")
   plf-grid (Tem)
   plt. show ()
  plt-figure (figsige = (10,5))
    sns. histplot ( of [ ' sulfase_year], bins = 30, kde = Palse)
   plt-HHe ("Distribution of Release Years")
   Plt. xlabell ("Release Years")
    plt- xlabel (" Nomber of Titles")
    ptshow ()
    Sne. countplot (data = df, x = 'type', palette = 'Set z')
   plt. Hitle ( " count of Hovies and Tu shows")
    plt. show()
     top- 100ntries = of ['lountry']. Value_ (ounts). head(10)
    top - countries . plot (tind = 'bas', colon = 'sky blw)
    plt-title ("Top 10 countries by Number of Titles")
   plt. ylabel 14 (cont")
    plt. Xtilbs (notation = 45)
    plt-show()
```

350-200 -150 100 50 2011 2013 2009 2015 2017 2019 Comment of the contract of the state of the Market of Shell a palace? The planning of the land The return Contraction of the party



WHITE PERSONNENT A SERVE - PA

(and martagon) - entity + Ha