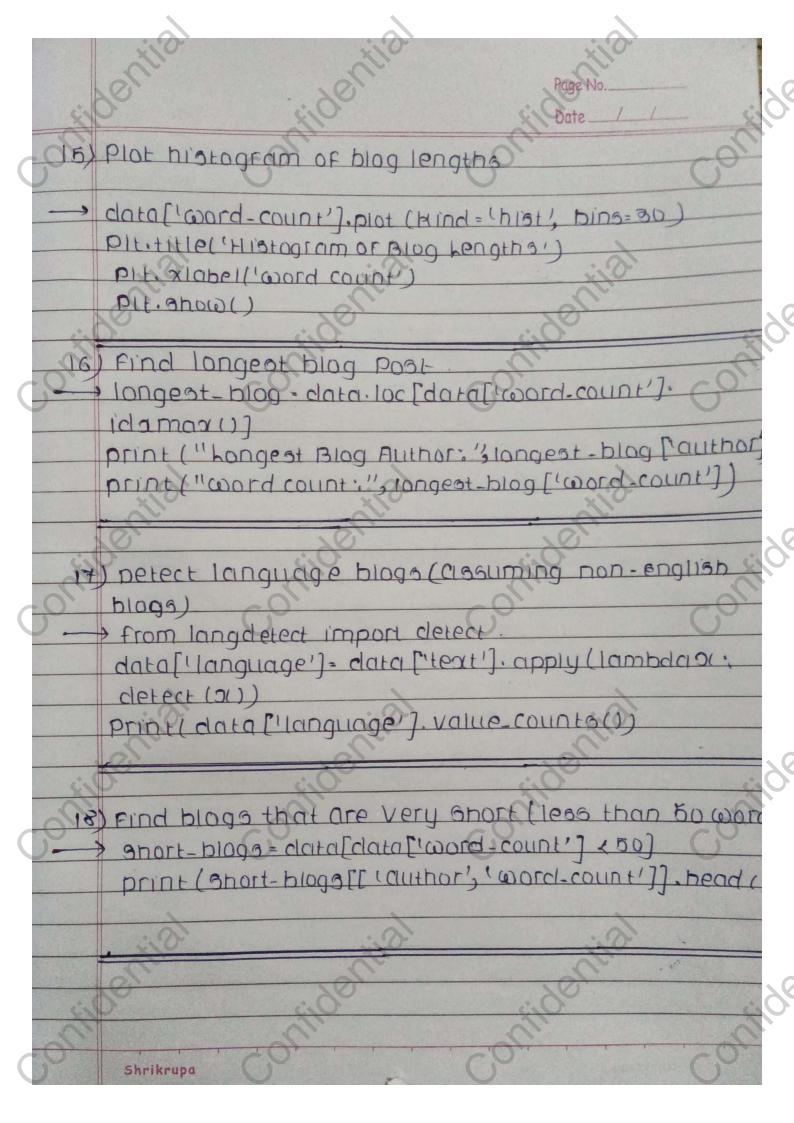
Name:- Arpita Nilean Tap Page No. PRN: - 202401070004 Date 98/04/2025 Batch- ET13 Topic: The Blog Authorship Corpus . Load the Blog Authorship corpus file import Pandas as Pd. # Assuming the dataset is in car format data = Pd. read_csv ('blog_authorship_corpus.c sv') print(data.head()) 2) Display total no of blogs print (f" Total number of blogs; slen(data)?") 3) Find all unique authors > unique_authors = data['author'], unique() print (P" Number of unique authors: Sien (unique. authors) 7" us Find number of blogs written by males male_blogs = data[data['gender'] = male'] print (P" number of blogs by males: & len (male blogs) Find average age of authors average - age = data['age/]. mean() print (f" Average age of authors : & average age : . of?") Shrikrupa

6) Find the most common industry among authors > common industry = data['industry].mode()[0] print (& most common industry: & common industry?") T) create a bar plot showing gender distribution import matplotlib pyplot as pit data['gender'] value counts(), plat (Hind='bar') PIT. title (Grender Distribution) PIL . Ollabel (Gender') PIt. ylabel ('name of Blogs') PIL. Show() 8) create a new country blug Length in words = [data[word conot] = dota['clean_teat'].apply (lambdax: len(x.split()); Print (datafficord-count')], describe()) 9) Find the author with maximum blog posts + top author = data l'author]. value-counts () print (f" Author with most blogs: § top author?") 10) Find average blog Length by gender > avg-length-by-gender=data.gouphy('gender') ['word-count']. mean() print (avg-length-by-gender) Shrikrupa

corre print eaun false print .csix la calcu unde unde print perce prin perce prin perce prin eau def i def i def i	correlation blw age lation = clata['age'].c leaned data to contended data data data data data data data d	orr(data['word-count']) on age and word. "')
11) Find - correption print coin eain false print coiv data false print coiv data print coiv coiv data print coiv coiv data print coiv coiv coiv data print coiv coiv coiv data print coiv coi	correlation blw age lation = clata['age'].c leaned data to contended data data data data data data data d	and blog kength arr(data['word-count']) an age and word. ")
Corre print eain 13 Calci print .c3V 13 Calci unde perce prin E per th Clea def i re cla	lation = clata['age'].c. (P''correlation between t: Scorrelation: . 2f3 cleaned data to co. to_cov('cleaned-blog	and blog kength orr(data['word count']) on age and word. "")
Corre print edita false print .c3v l3 calcu unde perce prin	lation = clata['age'].c. (P''correlation between t: Scorrelation: . 2f3 cleaned data to co. to_cov('cleaned-blog	orr(data['word count']) on age and word. "')
Print coin Print coin Print coiv Print coiv Print coiv Print coiv Print perce Prin Eper	cleaned data to co	n age and word.
eaun 12) Save data False print .c3v l3 / calcu unde unde perce prin perce perce prin perce	cleaned data to ca	
Print Calculate Print Calculate Print Linde Perce Print E per def i clea	cleaned data to co	V C O C O C
- data False print .csv 13 / calculate unde unde perce prin E per def i re data	· to_cov('cleaned-blog	-cuthorship.cov'sinder=
- data False print .csv lis calculate unde unde perce prin & perce p	· to_cov('cleaned-blog	-cuthorship.cov; index=
- data False print .csv 13 calculate unde unde perce prin & perce pr	· to_cov('cleaned-blog	-cluthorship. cov; inder=
False print .c3 unde unde perce prin perce perce prin perce prin perce p		
inde unde unde unde perce prin perce perce prin perce prin perce perce prin perce perce prin perce per		
inde unde unde unde perce prin perce	Kucleaned data saved	to ! cleaned - blog - authorship
Linde Perce prin perce perce prin perce perce prin perce perce perce prin perce perce perce prin perce pe		
unde unde unde perce prin def i re dai		
unde unde unde perce prin def i re dai		
perce prin perce perce prin perce perce prin perce prin perce perc	ulate percentage of blo	go written by authors
perce prin & per & per der c der c		C_{ρ}
prin & per iu clea der o re	er- 05 = data[data['age'	
th clear der control der contr	entage = (len(under-20)	
def of the def		is by authors under 25:
def (centage: of } 7."	
def (Laborationa and non
re	in text : remove specia	renarrans and nos
da	turn re oub (r/ \a-78-	7/9/11 text1
	filli le ano (1. fun-cu	ulitentij appinicious tour
PI	1 - [alana + anx 1] - dat	
	ta [clean - text] = dat	m- cext ff records)
	int (dataffitext); cle	
10)		
1		18/10
Shrike		i de l'ille
5	int (data[['text', 'cle	



Drivoential Page No. Date 19) Vialize blog counts by age group 3n3 countplot (x = 'age group') data = data, parette= 'sets' Pit title ('Blog counts by Age Group') OPIT. Show() 20) compute readability score (Flesch Reading 805e impart textstat data[readability'] = data[rean-text'] apply (textstat. flesch-reading-ease) print (data [['readability']].describe())