	1. Get the dataset
<pre>In [1]: Out[1]:</pre>	import pandas as pd df = pd.read_csv('papers.csv') df.head() id year title event_type pdf_name abstract paper_text Self-Organization of Associative Database and A Mean Field Theory A Mean Field Theory A Mean Field Theory 1-self-organization of of-associative database-an A Mean Field Theory 10-a-mean-field-Abstract Abstract Ab
	1 10 1987 of Layer IV of Visual Cort NaN theory-of-layer-iv-of-visual-c Storing Covariance by the Associative Long-Ter Bayesian Query Construction for Neural Network NaN theory-of-layer-iv-of-visual-c 100-storing-covariance-by-the-associative-long NaN covariance-by-the-associative-long 100-bayesian-query-construction-for-neural-ne NaN query-construction-for-neural-ne NaN theory-of-layer-iv-of-visual-c 100-storing-COVARIANCE BY Missing Nan theory-of-layer-iv-of-layer-iv-of-visual-c Nan covariance-by-the-associative-long Nan query-construction-for-neural-ne Nan theory-of-layer-iv-of-visual-c Nan covariance-by-the-associative-long Nan theory-of-layer-iv-of-layer-iv-of-visual-c Nan theory-of-layer-iv-of-visual-c Nan theory-of-layer-iv-of-visual-c Nan theory-of-layer-iv-of-visual-c Nan theory-of-layer-iv-of-visual-c Nan theory-of-layer-iv-of-visual-c Nan theory-of-layer-iv-of-visual-c Nan theory-of-visual-c Nan theory-of-visual-c Nan theory-of-visual-c Nan theory-of-visua
<pre>In [2]: Out[2]:</pre>	id year title event_type pdf_name abstract paper_text 7236 994 1994 Single Transistor Learning Synapses NaN transistor-learning-synapses.pdf Synapses\n\nPaul Ha
	Page 1994 Bias, Variance and the Combination of Least Sq Page 1994 Bias, Variance and the Combination of Least Sq Page 1994 A Real Time Clustering CMOS Neural Engine Page 1995 Bias, Variance and the Combination of Least Abstract Missing Combination of NaN Sing Combination of NaN Clustering Chos Neural Engine Clustering CMOS Neural Engine Chos Neural Engine Page 1998-learning Description of NaN Page 1998-learning Chos Nation Chass Page 1996 Bias, Variance Abstract Abstract Missing Chos Nation Clustering Chos Nation Chos Nation Chass A Real Time Clustering Chos Nation Clustering Chos Nation Chos Nation Chass Page 2996-bias-variance Abstract Missing Chos Nation Chast Chos Nation Chos Nation Chass A Real Time Clustering Chos Nation
In [3]:	Correlation and 999-correlation-and-interpolation Networks for Rea Man and-interpolation-and-interpolation-networks-for Missing Correlation and Interpolation Networks for Rea Networks for\nRe df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 7241 entries, 0 to 7240 Data columns (total 7 columns): # Column Non-Null Count Dtype</class>
	0 id 7241 non-null int64 1 year 7241 non-null int64 2 title 7241 non-null object 3 event_type 2422 non-null object 4 pdf_name 7241 non-null object 5 abstract 7241 non-null object 6 paper_text 7241 non-null object dtypes: int64(2), object(5) memory usage: 396.1+ KB
<pre>In [4]: Out[4]: In [5]: Out[5]:</pre>	df.shape (7241, 7) df.describe() id year count 7241.000000 7241.000000
	mean 3655.912167 2006.439718 std 2098.435219 8.759919 min 1.000000 1987.000000 25% 1849.000000 2000.000000 50% 3659.000000 2009.000000 75% 5473.000000 2014.000000 max 7284.000000 2017.000000
<pre>In [6]: Out[6]:</pre>	2. Dropping the unnecessary columns df.drop(['title', 'event_type', 'pdf_name', 'paper_text'], axis = 1, inp lace = True) df.head() id year abstract 1. 1097 Abstract Missing
	 1 1987 Abstract Missing 1 10 1987 Abstract Missing 2 100 1988 Abstract Missing 3 1000 1994 Abstract Missing 4 1001 1994 Abstract Missing 3. Dropping the rows where Abstract is missing
In [7]: Out[7]:	<pre>id year</pre>
In [8]: Out[8]:	2385 3164 2007 We present the first truly polynomial algorith 2388 3167 2007 Semi-supervised inductive learning concerns ho dataset.tail() id year abstract 6943 7280 2017 We revisit the classical analysis of generativ 6944 7281 2017 PAC maximum
In [9]: Out[9]:	6945 7282 2017 Community detection, which focuses on clusteri 6946 7283 2017 We propose a general framework for interactive 6947 7284 2017 We consider maximum likelihood estimation of l dataset.shape (3924, 3)
n [10]:	<pre>dataset['word_count'] = dataset['abstract'].apply(lambda x: len(str(x).s plit(" "))) dataset[['abstract', 'word_count']].head() C:\Users\kruti\anaconda3\lib\site-packages\ipykernel_launcher.py:1: Sett ingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value instead</pre>
ut[10]:	See the caveats in the documentation: https://pandas.pydata.org/pandas-d ocs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy """Entry point for launching an IPython kernel. abstract word_count 941 Non-negative matrix factorization (NMF) has pr 107 1067 Spike-triggered averaging techniques are effec 81 2384 It is known that determinining whether a DEC-P 67
n [11]: ut[11]:	2385 We present the first truly polynomial algorith 143 2388 Semi-supervised inductive learning concerns ho 119 dataset.word_count.describe() count 3924.000000 mean 148.390928 std 45.605755 min 19.000000 25% 116.000000
n [12]:	110.000000 50% 143.000000 75% 177.000000 max 317.000000 Name: word_count, dtype: float64 5. Finding the most common words frequency = pd.Series(' '.join(dataset['abstract']).split()).value_count s()[:20]
ut[12]:	the 29793 of 20918 a 16339 and 13626 to 12869 in 8980 that 7838 is 7666 for 7169
	We6238on5579we5167with4512this3677as3643are3529an3366by3197can2953
n [13]: ut[13]:	<pre>learning</pre>
:[د	bound-based 1 (GCG)-type 1 communities/blocks. 1 (STA), 1 non-metric, 1 justification: 1 non-additive 1 registered, 1 et 1 bels. 1 LSVMs 1
	<pre>1 confidence-set</pre>
n [14]:	7. Text Preprocessing 7. 1. Importing the necessary libraries for text preprocessing import re import nltk #nltk.download('stopwords') from nltk.corpus import stopwords from nltk.stem.porter import PorterStemmer
n [15]:	<pre>from nltk.tokenize import RegexpTokenizer #nltk.download('wordnet') from nltk.stem.wordnet import WordNetLemmatizer 7. 2. Creating a list of 'Custom' stop words ##Creating a list of stop words and adding custom stopwords stop_words = set(stopwords.words("english")) ##Creating a list of custom stopwords</pre>
n [18]:	<pre>new_words = ["using", "show", "result", "large", "also", "iv", "one", "t wo", "new", "previously", "shown"] stop_words = stop_words.union(new_words)</pre> 7.3. Cleaning the text
	<pre>text = dataset.iloc[i, 2].replace(r'[^\w\s]', ' ') #Convert to lowercase text = text.lower() #remove tags text=re.sub("</?.*?>"," <> ",text) # remove special characters and digits text=re.sub("(\\d \\\\))+"," ",text)</pre>
	<pre>##Convert to list from string text = text.split() ##Stemming ps=PorterStemmer() #Lemmatisation lem = WordNetLemmatizer() text = [lem.lemmatize(word) for word in text if not word in</pre>
	#View corpus item corpus[222] 'actor critic algorithm reinforcement learning achieving renewed popular ity due good convergence property situation approach often fail e g func tion approximation involved interestingly growing evidence actor critic approach based phasic dopamine signal play key role biological learning cortical basal ganglion derive temporal difference based actor critic le arning algorithm convergence proved without assuming separate time scale
n [21]:	actor critic approach demonstrated applying network spiking neuron estab lished relation phasic dopamine temporal difference signal lends support biological relevance algorithm' 8. Word Cloud #Word cloud from os import path from PIL import Image
	<pre>from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator import matplotlib.pyplot as plt %matplotlib inline wordcloud = WordCloud(</pre>
	fig = plt.figure(1) plt.imshow(wordcloud) plt.axis('off') plt.show() fig.savefig("word1.png", dpi=900) <wordcloud.wordcloud.wordcloud 0x000001879cdf4dc8="" at="" object=""> learning structure based framework tigstands state art feature</wordcloud.wordcloud.wordcloud>
	9. Text Preparation
n [22]:	 Tokenization - for this we use the Bag of Words model Vectorization - The list of words is then converted to a matrix of integers by the process of vectorisation. Vectorisation is also called feature extraction. We will use the CountVectoriser to tokenise the text and build a vocabulary of known words. from sklearn.feature_extraction.text import CountVectorizer import re cv=CountVectorizer(max_df=0.8, stop_words=stop_words, max_features=10000, ngram_range=(1,3))
n [23]: ut <mark>[23]</mark> :	<pre>X=cv.fit_transform(corpus) list(cv.vocabularykeys())[:10] ['non', 'negative', 'matrix', 'factorization', 'nmf', 'useful', 'decomposition',</pre>
n [25]:	<pre>'multivariate', 'data', 'different'] 10. Visualizing Uni-grams #Most frequently occuring words def get_top_n_words(corpus, n=None): vec = CountVectorizer().fit(corpus)</pre>
	<pre>bag_of_words = vec.transform(corpus) sum_words = bag_of_words.sum(axis=0) words_freq = [(word, sum_words[0, idx]) for word, idx in</pre>
ut[25]:	<pre>#Barplot of most freq words import seaborn as sns sns.set(rc={'figure.figsize':(13,8)}) g = sns.barplot(x="Word", y="Freq", data=top_df) g.set_xticklabels(g.get_xticklabels(), rotation=30) [Text(0, 0, 'model'), Text(0, 0, 'algorithm'), Text(0, 0, 'learning'), Text(0, 0, 'method'), Text(0, 0, 'problem'), Text(0, 0, 'data'),</pre>
	Text(0, 0, 'approach'), Text(0, 0, 'function'), Text(0, 0, 'based'), Text(0, 0, 'network'), Text(0, 0, 'time'), Text(0, 0, 'paper'), Text(0, 0, 'result'), Text(0, 0, 'task'), Text(0, 0, 'distribution'), Text(0, 0, 'propose'), Text(0, 0, 'state'),
	Text(0, 0, 'feature'), Text(0, 0, 'image'), Text(0, 0, 'performance')]
	2000
n [261·	11. Visualizing Bi-grams #Most frequently occuring Bi-grams def get top n2 words (corpus, n=None):
n [26]:	<pre>def get_top_n2_words(corpus, n=None): vec1 = CountVectorizer(ngram_range=(2,2),</pre>
	<pre>top2_df = pd.DataFrame(top2_words) top2_df.columns=["Bi-gram", "Freq"] print(top2_df) #Barplot of most freq Bi-grams import seaborn as sns sns.set(rc={'figure.figsize':(13,8)}) h=sns.barplot(x="Bi-gram", y="Freq", data=top2_df) h.set_xticklabels(h.get_xticklabels(), rotation=45) Bi-gram Freq state art 718</pre>
	neural network 561 machine learning 413 real world 376 learning algorithm 350 high dimensional 326 lower bound 259 data set 239 paper propose 237 optimization problem 225 graphical model 224 experimental result 220 loss function 208
ut[26]:	13 reinforcement learning 205 14 low rank 202 15 gradient descent 188 16 gaussian process 184 17 learning problem 180 18 convergence rate 178 19 stochastic gradient 175 [Text(0, 0, 'state art'), Text(0, 0, 'neural network'),
	<pre>Text(0, 0, 'machine learning'), Text(0, 0, 'real world'), Text(0, 0, 'learning algorithm'), Text(0, 0, 'high dimensional'), Text(0, 0, 'lower bound'), Text(0, 0, 'data set'), Text(0, 0, 'paper propose'), Text(0, 0, 'optimization problem'), Text(0, 0, 'graphical model'), Text(0, 0, 'experimental result'), Text(0, 0, 'loss function'),</pre>
	Text(0, 0, 'reinforcement learning'), Text(0, 0, 'low rank'), Text(0, 0, 'gradient descent'), Text(0, 0, 'gaussian process'), Text(0, 0, 'learning problem'), Text(0, 0, 'convergence rate'), Text(0, 0, 'stochastic gradient')]
	500 g 400 300
	100 State of the party of the
n [28]:	<pre>#Most frequently occuring Tri-grams def get_top_n3_words(corpus, n=None): vec1 = CountVectorizer(ngram_range=(3,3),</pre>
	<pre>words_freq = [(word, sum_words[0, idx]) for word, idx in</pre>
	<pre>sns.set(rc={'figure.figsize':(13,8)}) j=sns.barplot(x="Tri-gram", y="Freq", data=top3_df) j.set_xticklabels(j.get_xticklabels(), rotation=45) Tri-gram Freq 0 convolutional neural network 100 1 deep neural network 96 2 stochastic gradient descent 84 3 state art performance 82 4 real world datasets 76 5 recurrent neural network 73 6 state art method 73</pre>
	7 real world data 68 8 markov decision process 63 9 low rank matrix 63 10 synthetic real world 62 11 semi supervised learning 61 12 multi armed bandit 61 13 outperforms state art 60 14 state art result 60 15 latent variable model 57 16 markov chain monte 55 17 chain monte carlo 55
ut[28]:	18 principal component analysis 50 19 empirical risk minimization 50 [Text(0, 0, 'convolutional neural network'), Text(0, 0, 'deep neural network'), Text(0, 0, 'stochastic gradient descent'), Text(0, 0, 'state art performance'), Text(0, 0, 'real world datasets'), Text(0, 0, 'recurrent neural network'), Text(0, 0, 'state art method'), Text(0, 0, 'real world data'),
	<pre>Text(0, 0, 'markov decision process'), Text(0, 0, 'low rank matrix'), Text(0, 0, 'synthetic real world'), Text(0, 0, 'semi supervised learning'), Text(0, 0, 'multi armed bandit'), Text(0, 0, 'outperforms state art'), Text(0, 0, 'state art result'), Text(0, 0, 'latent variable model'), Text(0, 0, 'markov chain monte'), Text(0, 0, 'chain monte carlo'), Text(0, 0, 'principal component analysis'),</pre>
	40 - 20 - 20
	13. Using TF-IDF for improving the word_counts • TF — term frequency TF = (Frequency of a term in a document)/(total no. of terms in
n [29]:	<pre>document) • IDF — Inverse document frequency IDF = (log(Total documents))/(no. of documents with the term) from sklearn.feature_extraction.text import TfidfTransformer tfidf_transformer=TfidfTransformer(smooth_idf=True, use_idf=True) tfidf_transformer.fit(X) # get feature names</pre>
n [30]:	<pre>feature_names=cv.get_feature_names() # fetch document for which keywords needs to be extracted doc=corpus[532] #generate tf-idf for the given document tf_idf_vector=tfidf_transformer.transform(cv.transform([doc])) #Function for sorting tf_idf in descending order from scipy.sparse import coo_matrix def sort_coo(coo_matrix):</pre>
	<pre>def sort_coo(coo_matrix): tuples = zip(coo_matrix.col, coo_matrix.data) return sorted(tuples, key=lambda x: (x[1], x[0]), reverse=True) def extract_topn_from_vector(feature_names, sorted_items, topn=10): """get the feature names and tf-idf score of top n items""" #use only topn items from vector sorted_items = sorted_items[:topn] score_vals = [] feature_vals = []</pre>

score_vals.append(round(score, 3))
feature_vals.append(feature_names[idx])

results[feature_vals[idx]]=score_vals[idx]

keywords=extract_topn_from_vector(feature_names, sorted_items, 5)

Abstract:
present theory compositionality stochastic optimal control showing task optimal controller constructed certain primitive primitive feedback cont roller pursuing agenda mixed proportion much progress making towards age nda compatible agenda present task resulting composite control law prova bly optimal problem belongs certain class class rather general yet numbe r unique property bellman equation made linear even non linear discrete dynamic give rise compositionality developed special case linear dynamic gaussian noise framework yield analytical solution e non linear mixture linear quadratic regulator without requiring final cost quadratic genera lly natural set control primitive constructed applying svd green function bellman equation illustrate theory context human arm movement idea opt imality compositionality prominent field motor control yet hard reconcil e work make possible

#sort the tf-idf vectors by descending order of scores

#create a tuples of feature, score

return results

now print the results

print(k,keywords[k])

print("\nAbstract:")
print(doc)
print("\nKeywords:")

Keywords: compositionality 0.33 primitive 0.295

for k in keywords:

Abstract:

#results = zip(feature_vals, score_vals)

results= {}

for idx in range(len(feature_vals)):

sorted_items=sort_coo(tf_idf_vector.tocoo())

#extract only the top n; n here is 10