In [1]: **import** numpy **as** np import pandas as pd import os for dirname, \_, filenames in os.walk('/kaggle/input'): **for** filename **in** filenames: print(os.path.join(dirname, filename)) movies = pd.read\_csv('tmdb\_5000\_movies.csv') In [2]: credits = pd.read\_csv('tmdb\_5000\_credits.csv') movies.shape In [3]: (4803, 20)Out[3]: credits.head() In [4]: Out[4]: movie\_id title cast crew [{"credit\_id": [{"cast\_id": 242, "character": 0 19995 Avatar "52fe48009251416c750aca23", "Jake Sully", "... "de... [{"credit\_id": [{"cast\_id": 4, "character": Pirates of the Caribbean: 285 1 "52fe4232c3a36847f800b579", At World's End "Captain Jack Spa... [{"credit id": [{"cast id": 1, "character": 2 206647 Spectre "54805967c3a36829b5002c41", "James Bond", "cr... "de... [{"credit\_id": [{"cast\_id": 2, "character": "52fe4781c3a36847f81398c3", 3 49026 The Dark Knight Rises "Bruce Wayne / Ba... movies = movies.merge(credits, on='title') In [5]: movies.head() In [6]: budget homepage keywords original\_langua genres Out[6]: [{"id": [{"id": 28, 1463, 'name": "name": "Action"}, 237000000 http://www.avatarmovie.com/ 19995 "culture {"id": 12, clash"}, "nam... {"id":... [{"id": 270, [{"id": 12, "name": "name": "ocean"}, 300000000 285 http://disney.go.com/disneypictures/pirates/ "Adventure"}, {"id": 726, {"id": 14, "... "na... [{"id": 28, [{"id": 470, "name": "name": 245000000 "Action"}, http://www.sonypictures.com/movies/spectre/ 206647 "spy"}, {"id": 818, {"id": 12, "nam... "name... [{"id": 849, [{"id": 28, "name": "name": "dc "Action"}, 3 250000000 http://www.thedarkknightrises.com/ 49026 comics"}, {"id": 80. {"id": "nam... 853,... [{"id": 818, [{"id": 28, "name": "name": 260000000 "based on "Action"}, http://movies.disney.com/john-carter 49529 {"id": 12, novel"}, "nam... {"id":... 5 rows × 23 columns import ast In [7]: In [8]: def convert(text): L = []for i in ast.literal\_eval(text): L.append(i['name']) return L movies.dropna(inplace=True) In [9]: movies['genres'] = movies['genres'].apply(convert) In [10]: movies.head() id keywords original language budget homepage genres Out[10]: [{"id": [Action, 1463, Adventure, "name": 237000000 Fantasy, http://www.avatarmovie.com/ 19995 er "culture Science clash"}, Fiction] {"id":... [{"id": 270, [Adventure, "name": 300000000 "ocean"}, 285 http://disney.go.com/disneypictures/pirates/ Fantasy, eı {"id": 726, Action] "na... [{"id": 470, [Action, "name": 2 245000000 http://www.sonypictures.com/movies/spectre/ 206647 Adventure, "spy"}, {"id": 818, Crime] "name... [{"id": 849, [Action, "name": Crime, "dc 250000000 http://www.thedarkknightrises.com/ 49026 eı comics"}, Drama, Thriller] {"id": 853,... [{"id": 818, [Action. "name": Adventure, 260000000 49529 "based on http://movies.disney.com/john-carter er Science novel"}, Fiction] {"id":... 5 rows × 23 columns movies['keywords'] = movies['keywords'].apply(convert) In [49]: movies.head() Traceback (most recent call last) ValueError Input In [49], in <cell line: 1>() ----> 1 movies['keywords'] = movies['keywords'].apply(convert) 2 movies.head() File ~\anaconda3\lib\site-packages\pandas\core\series.py:4433, in Series.apply(se 1f, func, convert\_dtype, args, \*\*kwargs) 4323 **def** apply( 4324 self, 4325 func: AggFuncType,  $(\ldots)$ 4328 \*\*kwargs, 4329 ) -> DataFrame | Series: 4330 4331 Invoke function on values of Series. 4332  $(\ldots)$ 4431 dtype: float64 4432 return SeriesApply(self, func, convert\_dtype, args, kwargs).apply() -> 4433 File ~\anaconda3\lib\site-packages\pandas\core\apply.py:1082, in SeriesApply.appl 1078 if isinstance(self.f, str): 1079 # if we are a string, try to dispatch 1080 return self.apply\_str() -> 1082 return self.apply\_standard() File ~\anaconda3\lib\site-packages\pandas\core\apply.py:1137, in SeriesApply.appl y\_standard(self) 1131 values = obj.astype(object).\_values 1132 # error: Argument 2 to "map\_infer" has incompatible type # "Union[Callable[..., Any], str, List[Union[Callable[..., Any], 1133 str]], # Dict[Hashable, Union[Union[Callable[..., Any], str], 1134 # List[Union[Callable[..., Any], str]]]]]"; expected 1135 # "Callable[[Any], Any]" 1136 mapped = lib.map\_infer( -> 1137 values, 1138 1139 f, # type: ignore[arg-type] 1140 convert=self.convert\_dtype, 1143 **if** len(mapped) **and** isinstance(mapped[0], ABCSeries): 1144 # GH#43986 Need to do list(mapped) in order to get treated as nested # See also GH#25959 regarding EA support 1145 1146 return obj.\_constructor\_expanddim(list(mapped), index=obj.index) File ~\anaconda3\lib\site-packages\pandas\\_libs\lib.pyx:2870, in pandas.\_libs.li b.map\_infer() Input In [45], in convert(text) 1 def convert(text): 2 L = []for i in ast.literal\_eval(text): ---> 3 4 L.append(i['name']) 5 return L File ~\anaconda3\lib\ast.py:105, in literal\_eval(node\_or\_string) 103 return left - right 104 return \_convert\_signed\_num(node) --> 105 return \_convert(node\_or\_string) File ~\anaconda3\lib\ast.py:104, in literal\_eval.<locals>.\_convert(node) 102 else: 103 **return** left - right --> 104 return \_convert\_signed\_num(node) File ~\anaconda3\lib\ast.py:78, in literal\_eval.<locals>.\_convert\_signed\_num(nod e) 76 else: 77 **return** - operand ---> 78 return \_convert\_num(node) File ~\anaconda3\lib\ast.py:69, in literal\_eval.<locals>.\_convert\_num(node) 67 def \_convert\_num(node): if not isinstance(node, Constant) or type(node.value) not in (int, fl oat, complex): ---> 69 \_raise\_malformed\_node(node) return node.value File ~\anaconda3\lib\ast.py:66, in literal\_eval.<locals>.\_raise\_malformed\_node(no 65 def \_raise\_malformed\_node(node): raise ValueError(f'malformed node or string: {node!r}') ValueError: malformed node or string: ['culture clash', 'future', 'space war', 's
pace colony', 'society', 'space travel', 'futuristic', 'romance', 'space', 'alie
n', 'tribe', 'alien planet', 'cgi', 'marine', 'soldier', 'battle', 'love affair', import ast ast.literal\_eval('[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, [{'id': 28, 'name': 'Action'}, Out[13]: {'id': 12, 'name': 'Adventure'}, {'id': 14, 'name': 'Fantasy'}, ['id': 878, 'name': 'Science Fiction'}] def convert3(text): In [14]: L = []counter = 0for i in ast.literal\_eval(text): if counter < 3:</pre> L.append(i['name']) counter+=1 return L In [15]: movies['cast'] = movies['cast'].apply(convert) movies.head() budget homepage keywords original\_language Out[15]: genres [culture [Action, clash, Adventure, future, 237000000 Fantasy, 19995 space http://www.avatarmovie.com/ er Science war, Fiction] space colon... [ocean, drug abuse, [Adventure. 1 300000000 Fantasy, http://disney.go.com/disneypictures/pirates/ 285 exotic Action] island. east india [spy, based on [Action, novel, **2** 245000000 Adventure, http://www.sonypictures.com/movies/spectre/ 206647 secret er Crime] agent, seauel. mi... [dc [Action, comics, Crime, crime 250000000 49026 http://www.thedarkknightrises.com/ er Drama, fighter, Thriller] terrorist. secret i... [based on novel. [Action, Adventure, mars, 260000000 http://movies.disney.com/john-carter 49529 er Science medallion, Fiction] space travel... 5 rows × 23 columns movies['cast'] = movies['cast'].apply(lambda x:x[0:3]) In [16]: In [17]: def fetch\_director(text): for i in ast.literal\_eval(text): if i['job'] == 'Director': L.append(i['name']) return L movies['crew'] = movies['crew'].apply(fetch\_director) In [18]: In [19]: def collapse(L): L1 = []for i in L: L1.append(i.replace(" ","")) return L1 movies['cast'] = movies['cast'].apply(collapse) In [20]: movies['crew'] = movies['crew'].apply(collapse) movies['genres'] = movies['genres'].apply(collapse) movies['keywords'] = movies['keywords'].apply(collapse) In [21]: movies.head() keywords original\_l Out[21]: budget homepage id genres [cultureclash, [Action, future, Adventure, 237000000 http://www.avatarmovie.com/ 19995 spacewar, Fantasy, spacecolony, ScienceFiction] [ocean, [Adventure, drugabuse, 300000000 Fantasy, http://disney.go.com/disneypictures/pirates/ 285 exoticisland, Action] eastindiatrad... [spy, [Action, basedonnovel, http://www.sonypictures.com/movies/spectre/ 206647 2 245000000 Adventure, secretagent, Crime] sequel, mi6, ... [dccomics, [Action, Crime, crimefighter, 250000000 http://www.thedarkknightrises.com/ 49026 Drama, terrorist, Thriller] secretiden... [basedonnovel, [Action, mars, Adventure, medallion, 260000000 http://movies.disney.com/john-carter 49529 ScienceFiction] spacetravel, p... 5 rows × 23 columns movies['overview'] = movies['overview'].apply(lambda x:x.split()) In [22]: movies['tags'] = movies['overview'] + movies['genres'] + movies['keywords'] + movi In [23]: new = movies.drop(columns=['overview', 'genres', 'keywords', 'cast', 'crew']) In [24]: In [25]: new['tags'] = new['tags'].apply(lambda x: " ".join(x)) new.head() id original\_language original\_title popula Out[25]: budget homepage 237000000 http://www.avatarmovie.com/ 19995 Avatar 150.437 en Pirates of the Caribbean: 300000000 285 139.082 http://disney.go.com/disneypictures/pirates/ At World's End 2 245000000 http://www.sonypictures.com/movies/spectre/ 206647 en Spectre 107.376 The Dark 3 250000000 http://www.thedarkknightrises.com/ 49026 112.312 en **Knight Rises** 4 260000000 43.926 http://movies.disney.com/john-carter 49529 John Carter en In [26]: **from** sklearn.feature\_extraction.text **import** CountVectorizer cv = CountVectorizer(max\_features=5000, stop\_words='english') vector = cv.fit\_transform(new['tags']).toarray() In [27]: vector.shape In [28]: (1494, 5000) Out[28]: In [29]: from sklearn.metrics.pairwise import cosine\_similarity In [30]: similarity = cosine\_similarity(vector) In [31]: similarity 0.08134892, 0.05423261, ..., 0. , 0.05504819, array([[1. Out[31]: 0.02469324], , 0.05882353, ..., 0.04428074, 0. [0.08134892, 1. [0.05423261, 0.05882353, 1. , ..., 0. 0. ], [0. , 0.04428074, 0. , 0. 0.04032389], , 0. [0.05504819, 0. , ..., 0. 0. ], , 0. , ..., 0.04032389, 0. [0.02469324, 0. ]]) In [32]: new[new['title'] == 'The Lego Movie'].index[0] Out[32]: def recommend(movie): In [33]: index = new[new['title'] == movie].index[0] distances = sorted(list(enumerate(similarity[index])), reverse=True, key = lambo for i in distances[1:6]: print(new.iloc[i[0]].title) recommend('Gandhi') In [37]: IndexError Traceback (most recent call last) Input In [37], in <cell line: 1>() ----> 1 recommend('Gandhi') Input In [33], in recommend(movie) 1 def recommend(movie): index = new[new['title'] == movie].index[0] 3 distances = sorted(list(enumerate(similarity[index])), reverse=True, ke  $y = lambda \times : \times [1]$ for i in distances[1:6]: File ~\anaconda3\lib\site-packages\pandas\core\indexes\base.py:5039, in Index.\_\_g etitem\_\_(self, key) 5036 if is\_integer(key) or is\_float(key): # GH#44051 exclude bool, which would return a 2d ndarray 5038 key = com.cast\_scalar\_indexer(key, warn\_float=True) -> 5039 return getitem(key) 5041 **if** isinstance(key, slice): # This case is separated from the conditional above to avoid 5043 # pessimization com.is\_bool\_indexer and ndim checks. 5044 result = getitem(key) IndexError: index 0 is out of bounds for axis 0 with size 0 In [34]: import pickle In [35]: pickle.dump(new,open('movie\_list.pkl','wb')) pickle.dump(similarity,open('similarity.pkl','wb')) recommend('Gandhi') In [36]: Traceback (most recent call last) IndexError Input In [36], in <cell line: 1>() ----> 1 recommend('Gandhi') Input In [33], in recommend(movie) 1 def recommend(movie): index = new[new['title'] == movie].index[0] 3 distances = sorted(list(enumerate(similarity[index])), reverse=True, ke y = lambda x: x[1])for i in distances[1:6]: File ~\anaconda3\lib\site-packages\pandas\core\indexes\base.py:5039, in Index.\_\_g etitem\_\_(self, key) 5036 if is\_integer(key) or is\_float(key): 5037 # GH#44051 exclude bool, which would return a 2d ndarray 5038 key = com.cast\_scalar\_indexer(key, warn\_float=True) return getitem(key) 5041 **if** isinstance(key, slice): # This case is separated from the conditional above to avoid 5043 # pessimization com.is\_bool\_indexer and ndim checks. result = getitem(key) 5044 IndexError: index 0 is out of bounds for axis 0 with size 0 In [ ]: