# Day 37

Jan 6

Explaining my understanding of
the Movie Recommendation Engine
by hand

1. Step : Importing neccessaries
libraries and dependencies

ML will be done using sci-rcit
learning. The libraries of interest are
→ Count vectorizer : change text.
documents to integers/vectors. Must
of the libraries work better with numbers
than text

→ Cosine similarity: find the similarity between files of interest

Also we will import pandas $\xi$ summary to work with dataframes. and data

→ Step 2: Import and read the data.

data = pd.read_csv("data.csv")
                                    → extension
↑                         ↓
Python command            └ data to read

.• Get details of file

data.info() :→ Get information about column headings of data, n. of data that are non-null.

Fill the null row's [ columns
data.columns() : Get info about the headings.

Step 3: Select the features to be used
for the                    inputs of interest   ' because test
                               ↓     ↑   ↓       they are test
→ features = [ ' a ' , ' b ' , ' e ' ]
                    ℗
              they are list

Step 4: Put all the features of interest in
   one    column                   ✓ converting nan to
                                              strings
         for  feature  in features:
         df [ feature ] = df [ feature ] . fillna (')

This is how this works.
              first feature is    a
              for    a    in features:
                 df [ a ] = df [ a ]. fillna (' ")
                          └─┘          └────────┘
                     in              Save this

We used an empty string .fillna(' ')
because the datatype is a string. If
numerical value, the appropriate type
will be 0

For definition

placeholder
↓
def combine_feature (row):
try:
return row ['a'] + " " + return row 'b'
↑
column name

return the
row of the
column name a +
" row of the
column name [b]

excepts
print ("Error:", row)

Creating column in data frame:

df [' combined feature] : df.apply (combine _feature,
axis: r )
↓
apply +
row

To access a particular data in the data frame.

$$df.iloc[0]['combined features']$$

dataframe function

row element & a(les)

column name

---

**Step 4: Convert text document to integer/ vectors**

⤷ initializing it

cv = CountVectorizer(-)

count_matrix = cv.fit_transform(df["combined_features"])

⤷ view in human readable form

a = count.matrix.toarray()

⤷ allows to view

Access vocabularies learnt by CountVectorizer

$$Vocab = CV.\ vocabulary\_.keys()$$

Cosine_similarity

vector of
combined feature

cos_sim = cosine_similarity(count_matrix)