# Phase 2: Innovation

## Introduction to Fake News Detection using NLP:

*Fake news, in today's information age, presents a significant challenge to the dissemination of accurate and trustworthy information. The deliberate spread of false or misleading information through various media channels can have far-reaching consequences, from influencing public opinion to undermining trust in journalism and institutions. To combat this issue, Natural Language Processing (NLP) techniques and technologies have emerged as valuable tools in identifying and mitigating the impact of fake news.*



*Fake news detection using NLP involves the application of advanced computational linguistic methods to distinguish between genuine, fact-based news and fabricated, misleading content. This is achieved by analyzing the textual content, context, and linguistic patterns within news articles, social media posts, and other textual sources. The goal is to develop automated systems that can accurately classify information as either real or fake, aiding both individuals and organizations in making informed decisions and preventing the spread of misinformation.*

## 1. Advanced Techniques and BERT Integration:

◆ *In this phase, you move beyond traditional NLP techniques and embrace advanced methods, with a primary focus on integrating BERT (Bidirectional Encoder Representations from Transformers). BERT is a state-of-the-art deep learning model in NLP known for its ability to capture complex relationships in text data.*

## 2. Tokenization:

◆ *BERT requires tokenization of input text. Tokenization is the process of breaking down a piece of text into smaller units, typically words or subwords. BERT handles variable-length sequences efficiently, and you'll need to integrate a BERT tokenizer into your project. BERT tokenization is subword-based and can capture fine-grained linguistic information.*

## 3. Model Architecture:

◆ Designing a BERT-based model architecture is a crucial step. BERT models are known for their deep and complex architectures. BERT's architecture includes multiple layers of attention mechanisms, which allows it to understand the context of words in a sentence, including the relationships between them. Your model should be designed for binary classification, distinguishing between real and fake news.

## 4. Training and Evaluation:

◆ Train the BERT-based model on your preprocessed data. BERT models can be computationally intensive due to their complexity. Training will likely require substantial computational resources. After training, evaluate the model's performance using appropriate metrics like accuracy, precision, recall, F1-score, and possibly the AUC-ROC curve to assess its ability to discriminate between real and fake news.

## 5. Hyperparameter Tuning:

◆ **Experiment with hyperparameter tuning to optimize the BERT-based model's performance. This may involve adjusting learning rates, batch sizes, and other hyperparameters specific to BERT. Fine-tuning BERT for your specific fake news detection task is essential for achieving the best results.**

## 6. Error Analysis:

◆ **Analyze the model's misclassifications to gain insights into areas where it struggles. This could involve examining the types of false positives and false negatives the model produces. Such an analysis can guide further improvements in your model or data preprocessing steps, helping you refine the model's performance.**

## 7. Deployment (if applicable):

◆ **Consider how you'll deploy the BERT-based model in a real-world application. Deployment can involve creating a web service or API that allows users to input news articles, and the model will make predictions on whether they are real or fake. Deploying deep learning models like BERT typically requires infrastructure to handle inference requests efficiently.**
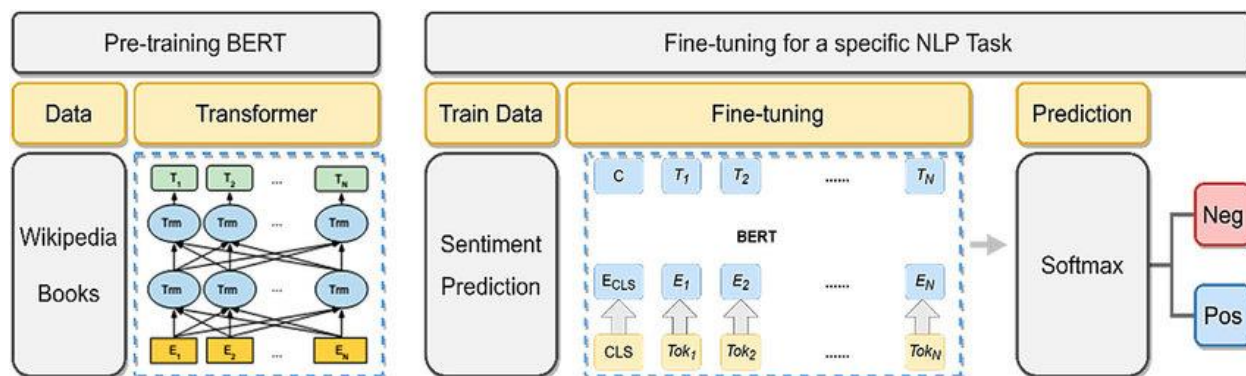
## Some Advanced Techniques

*BERT (Bidirectional Encoder Representations from Transformers):*

*BERT is a pre-trained deep learning model introduced by Google in 2018. It is designed to understand the context and semantics of words in a sentence by considering both left and right context simultaneously. This bidirectional understanding of text is a significant departure from earlier models that typically focused on either left-to-right or right-to-left language modeling.*

*Key features of BERT:*

◆ *Bidirectional Context*

◆ *Transformer Architecture*

◆ *Pre-training and Fine-Tuning*
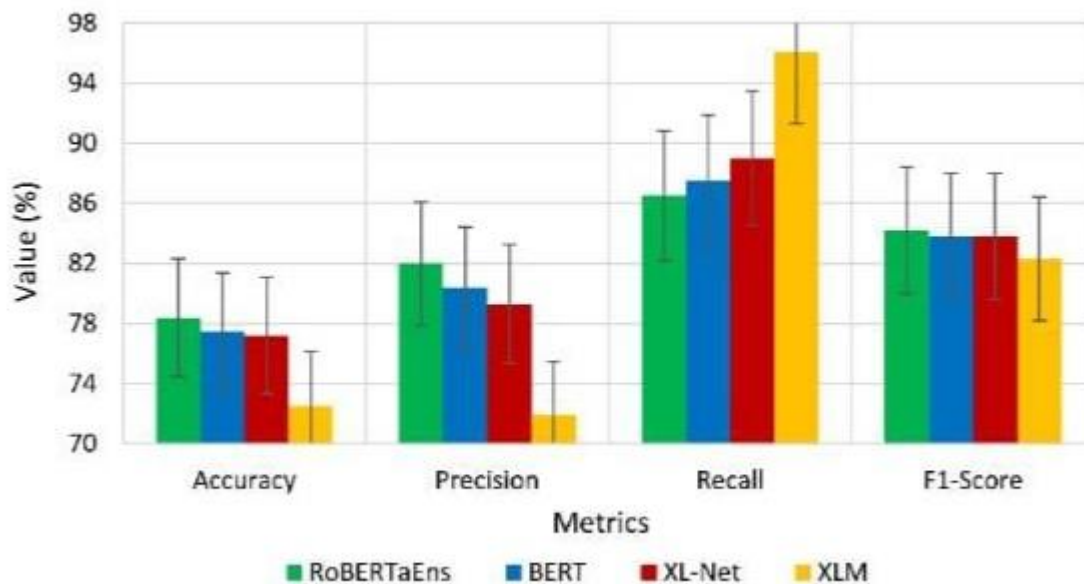
◆ *State-of-the-Art Results*

*RoBERTa (A Robustly Optimized BERT Pretraining Approach):*

*RoBERTa, introduced by Facebook AI in 2019, can be thought of as an extension and optimization of the BERT model. It aims to improve upon BERT's performance by addressing some of its limitations.*

*Key features of RoBERTa:*

◆ *Large-Scale Training*

◆ *Dynamically Masked Data*

◆ *Training Variant*

◆ *Improved Results*

# Program :

```
[16]:  import numpy as
       np import pandas
       as pd
       from nltk.corpus import stopwords
       import plotly.express as px
       from wordcloud import
       WordCloud
       from matplotlib import pyplot as plt
       from sklearn.model_selection import train_test_split
       from sklearn.feature_extraction.text import CountVectorizer
       import tensorflow as tf
       from tensorflow.keras.callbacks import EarlyStopping
       from sklearn.metrics import accuracy_score
```

```
[17]:  fake_news = pd.read_csv('Fake.csv')
       true_news = pd.read_csv('True.csv')
       fake_news.head(10)
```

```
[17]:                                                  title  \
       0    Donald Trump Sends Out Embarrassing New Year'...
```

```
1    Drunk Bragging Trump Staffer Started Russian ...
2    Sheriff David Clarke Becomes An Internet Joke...
3    Trump Is So Obsessed He Even Has Obama's Name...
4    Pope Francis Just Called Out Donald Trump Dur...
5    Racist Alabama Cops Brutalize Black Boy While...
6    Fresh Off The Golf Course, Trump Lashes Out A...
7    Trump Said Some INSANELY Racist Stuff Inside ...
8    Former CIA Director Slams Trump Over UN Bully...
9    WATCH: Brand-New Pro-Trump Ad Features So Muc...


                                          text subject  \
0  Donald Trump just couldn t wish all Americans ...    News
1  House Intelligence Committee Chairman Devin Nu...    News
2  On Friday, it was revealed that former Milwauk...    News
3  On Christmas day, Donald Trump announced that ...    News
4  Pope Francis used his annual Christmas Day mes...    News
```

```
5   The  number  of  cases  of  cops  brutalizing  and  ki...    News
6   Donald  Trump  spent  a  good  portion  of  his  day  a...    News
7   In  the  wake  of  yet  another  court  decision  that...    News
8   Many  people  have  raised  the  alarm  regarding  th...    News
9   Just  when  you  might  have  thought  we  d  get  a  br...    News


                    date
0   December  31,  2017
1   December  31,  2017
2   December  30,  2017
3   December  29,  2017
4   December  25,  2017
5   December  25,  2017
6   December  23,  2017
7   December  23,  2017
8   December  22,  2017
9   December  21,  2017
```

[18]: `true_news.head(10)`

[18]:
```
                                                title  \
0   As U.S. budget fight looms, Republicans flip t...
1   U.S. military to accept transgender recruits o...
2   Senior U.S. Republican senator: 'Let Mr. Muell...
3   FBI Russia probe helped by Australian diplomat...
4   Trump wants Postal Service to charge 'much mor...
5   White House, Congress prepare for talks on spe...
6   Trump says Russia probe will be fair, but time...
7   Factbox: Trump on Twitter (Dec 29) - Approval ...
8         Trump on Twitter (Dec 28) - Global Warming
9   Alabama official to certify Senator-elect Jone...


                                                text         subject  \
0   WASHINGTON (Reuters) - The head of a conservat...   politicsNews
1   WASHINGTON (Reuters) - Transgender people will...   politicsNews
2   WASHINGTON (Reuters) - The special counsel inv...   politicsNews
3   WASHINGTON (Reuters) - Trump campaign adviser ...   politicsNews
4   SEATTLE/WASHINGTON (Reuters) - President Donal...   politicsNews
5   WEST PALM BEACH, Fla./WASHINGTON (Reuters) - T...   politicsNews
6   WEST PALM BEACH, Fla (Reuters) - President Don...   politicsNews
7   The  following  statements  were  posted  to  the  ve...   politicsNews
8   The  following  statements  were  posted  to  the  ve...   politicsNews
9   WASHINGTON (Reuters) - Alabama Secretary of St...   politicsNews


                    date
0   December  31,  2017
1   December  29,  2017
```

```
2    December 31, 2017
3    December 30, 2017
4    December 29, 2017
5    December 29, 2017
6    December 29, 2017
7    December 29, 2017
8    December 29, 2017
9    December 28, 2017
```

[19] : `fake_news.info()`

```
<class    'pandas.core.frame.DataFrame'>
RangeIndex: 23481 entries, 0 to 23480
Data columns (total 4 columns):
 #    Column     Non-Null Count  Dtype
---   -------    --------------  ------
 0    title      23481 non-null  object
 1    text       23481 non-null  object
 2    subject    23481 non-null  object
 3    date       23481 non-null  object
dtypes: object(4)
memory  usage: 733.9+ KB
```

[20] : 
```
true_news['True']  =  1
fake_news['True']  =  0
```

[21] : `true_news.drop(columns=['title','subject','date'])`

[21]:

|       | text | True |
|-------|------|------|
| 0     | WASHINGTON (Reuters) - The head of a conservat... | 1 |
| 1     | WASHINGTON (Reuters) - Transgender people will... | 1 |
| 2     | WASHINGTON (Reuters) - The special counsel inv... | 1 |
| 3     | WASHINGTON (Reuters) - Trump campaign adviser ... | 1 |
| 4     | SEATTLE/WASHINGTON (Reuters) - President Donal... | 1 |
| ...   | ... | ... |
| 21412 | BRUSSELS (Reuters) - NATO allies on Tuesday we... | 1 |
| 21413 | LONDON (Reuters) - LexisNexis, a provider of l... | 1 |
| 21414 | MINSK (Reuters) - In the shadow of disused Sov... | 1 |
| 21415 | MOSCOW (Reuters) - Vatican Secretary of State ... | 1 |
| 21416 | JAKARTA (Reuters) - Indonesia will buy 11 Sukh... | 1 |

[21417 rows x 2 columns]

[22] : `fake_news.drop(columns=['title','subject','date'])`

[22] :

|   | text | True |
|---|------|------|
| 0 | Donald Trump just couldn t wish all Americans ... | 0 |

```
1        House Intelligence Committee Chairman Devin Nu...    0
2        On Friday, it was revealed that former Milwauk...    0
3        On Christmas day, Donald Trump announced that ...    0
4        Pope Francis used his annual Christmas Day mes...    0
...                                                      ...   ...
23476    21st Century Wire says As 21WIRE reported earl...    0
23477    21st Century Wire says It s a familiar theme. ...    0
23478    Patrick  Henningsen  21st Century WireRemember ...   0
23479    21st Century Wire says Al Jazeera America will...    0
23480    21st Century Wire says As 21WIRE predicted in ...    0

[23481 rows x 2 columns]
```

[23] : 
```
dataset = pd.concat([true_news, fake_news], axis=0)
clean_data = dataset.drop(columns=['title','subject','date'])
clean_data
```

[23]:
```
                                                    text    True
0        WASHINGTON (Reuters) - The head of a conservat...    1
1        WASHINGTON (Reuters) - Transgender people will...    1
2        WASHINGTON (Reuters) - The special counsel inv...    1
3        WASHINGTON (Reuters) - Trump campaign adviser ...    1
4        SEATTLE/WASHINGTON (Reuters) - President Donal...    1
...                                                      ...   ...
23476    21st Century Wire says As 21WIRE reported earl...    0
23477    21st Century Wire says It s a familiar theme. ...    0
23478    Patrick  Henningsen  21st Century WireRemember ...   0
23479    21st Century Wire says Al Jazeera America will...    0
23480    21st Century Wire says As 21WIRE predicted in ...    0

[44898 rows x 2 columns]
```

[24] :
```
clean_data.dtypes
```

[24] :
```
text     object
True      int64
dtype: object
```

[25] :
```
sub = dataset.groupby('subject').count()['title']
print(sub)
plt.figure(figsize=(10,10))
px.pie(dataset['subject'],names=dataset['subject'],title='Subject')
```

```
subject
Government News     1570
Middle-east          778
News                9050
```

1

```
US_News              783
left-news           4459
politics            6841
politicsNews       11272
worldnews          10145
Name: title,  dtype: int64
```

```
<Figure size  1000x1000 with 0 Axes>
```

[26] :
```
x =
clean_data.iloc[:,0]  y
= clean_data['True']
```

```
x :
 0     WASHINGTON (Reuters) - The head of a conservat...
 1     WASHINGTON (Reuters) - Transgender people will...
 2     WASHINGTON (Reuters) - The special counsel inv...
 3     WASHINGTON (Reuters) - Trump campaign adviser ...
 4     SEATTLE/WASHINGTON (Reuters) - President Donal...
 5     WEST PALM BEACH, Fla./WASHINGTON (Reuters) - T...
 6     WEST PALM BEACH, Fla (Reuters) - President Don...
 7     The following statements were posted to the ve...
 8     The following statements were posted to the ve...
 9     WASHINGTON (Reuters) - Alabama Secretary of St...
Name: text, dtype: object
 y :
 0    1
 1    1
 2    1
 3    1
 4    1
 5    1
 6    1
 7    1
 8    1
 9    1
Name: True, dtype: int64
```

[27] :
```
para = x.tolist()
words = "
".join(para)
```

[28] :
```
wordgroup = "".join(chars)
wordgroup[0:140]
```

[28] : 'WASHINGTON Reuters The head of a conservative Republican faction in the US Congress who voted this month for a huge expansion of the nation'

```
[29]:  plt.figure(figsize=(10,10))
        plt.imshow(WordCloud().generate(wordgroup))
```

[29]:  <matplotlib.image.AxesImage at 0x1f20cf0b070>



```
[30]:  print('number of words : ',len([word for word in wordgroup.split()]))
```

number of words :    18140003

```
[31]:  wordgroup.split()[0:10]
```

[31]:  ['WASHINGTON',
        'Reuters',
        'The',
        'head',
        'of',
        'a',
        'conservative',
        'Republican',
        'faction',
        'in']

```
[32]:  samp =
        clean_data.sample(n=3000) samp
```

[32]:                                               text    True
        9533    NEW YORK (Reuters) - Democratic Party activist...       1
        2131    If you ask any conservative, Ferguson police o...       0

```

```
21708   SPOT ON RACHEL ZONATION The Obamas pride thems...        0
2144    WASHINGTON (Reuters) - A U.S. congressional pa...        1
7178    If you thought Trump supporters wered bad, wai...        0
...                                                      ...     ...
15089   ERBIL, Iraq (Reuters) - Iraqi forces launched ...        1
13085   Besides Trump, no one has a bigger target on h...        0
11346   NEW DELHI, ISLAMABAD (Reuters) - India denounc...        1
4242    SAN FRANCISCO (Reuters) - A California federal...        1
8152    (Reuters) - Members of Maná, the Spanish-langu...        1

[3000 rows x 2 columns]
```

[33] :
```
truth_dist = samp.groupby('True').count()
truth_dist
```

[33] :
```
        text
True
0       1627
1       1373
```

[34] :
```
para_samp =
samp.iloc[:,0].tolist() group ="
".join(para_samp)
chars = [char for char in group.split() if char not in string.punctuation]
print('Number of words in this 3000 entry sample data : ',len(" ".join(chars).
```

Number of words in this 3000 entry sample data :    1207736

[35] :
```
word_samp  =   "   ".join(chars).split()
words = [word.lower() for word in word_samp]
words[0:20]
```

[35] :
```
['new',
 'york',
 '(reuters)',
 'democratic',
 'party',
 'activists',
 'in',
 'some',
 'u.s.',
 'states',
 'are',
 'using',
 'donald',
 'trump,',
 'the',
```

```
    'republican',
    'presidential',
    'candidate',
    'who',
    'has']
```

[36] : `len(words)`

[36]: 1207736

[37] : `samp.dtypes`

[37]: text       object
      True        int64
      dtype: object

[39]: `import nltk`

[40]: `nltk.download('stopwords')`

[nltk_data] Downloading package stopwords to
[nltk_data]        C:\Users\Administrator\AppData\Roaming\nltk_data...
[nltk_data]        Unzipping corpora\stopwords.zip.

[40]: True

[41]:
```python
import nltk
nltk.download('stopwords')

from nltk.corpus import stopwords

# Assuming you already have your 'words' list
imp_word = [word.lower() for word in words if word not in stopwords.
  words('english')]
imp_word[0:20]
```

[nltk_data] Downloading package stopwords to
[nltk_data]        C:\Users\Administrator\AppData\Roaming\nltk_data...
[nltk_data]      Package stopwords is already up-to-date!

[41]: ['new',
       'york',
       '(reuters)',
       'democratic',
       'party',
       'activists',
       'u.s.',

```
            'states',
            'using',
            'donald',
            'trump,',
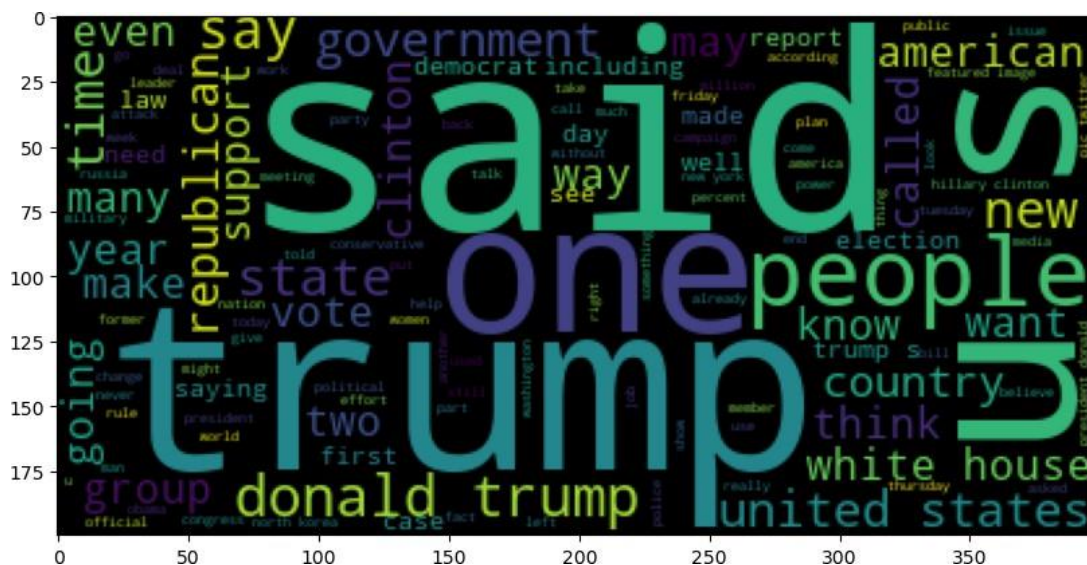            'republican',
            'presidential',
            'candidate',
            'stirred',
            'controversy',
            'comments',
            'illegal',
            'immigrants',
            'women,']
```

[42]:
```
imp_word =[word.lower() for word in words if word not in stopwords.
    words('english')]
imp_word[0:20]
```

[42]:
```
['new',
    'york',
    '(reuters)',
    'democratic',
    'party',
    'activists',
    'u.s.',
    'states',
    'using',
    'donald',
    'trump,',
    'republican',
    'presidential',
    'candidate',
    'stirred',
    'controversy',
    'comments',
    'illegal',
    'immigrants',
    'women,']
```

[43]:
```
plt.figure(figsize=(10,10))
plt.imshow(WordCloud().generate(" ".join(imp_word)))
```

[43] : <matplotlib.image.AxesImage at 0x1f2a5efc460>

```
[44]: vect =
CountVectorizer().fit_transform(para_samp).toarray() vect
```

```
[44]:  array([[0, 0,  0,  ..., 0, 0, 0],
           [0, 0,  0,  ..., 0, 0, 0],
           [0, 1,  0,  ..., 0, 0, 0],
           ...,
           [0, 0,  0,  ..., 0, 0, 0],
           [0, 0,  0,  ..., 0, 0, 0],
           [0, 0,  0,  ..., 0, 0, 0]],  dtype=int64)
```

```
[45]: vect_data =
pd.DataFrame(vect) vect_data
```

[45]:

|      | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | \ |
|------|---|---|---|---|---|---|---|---|---|---|---|
| 0    | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 1    | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2    | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 3    | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 4    | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| ...  | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 2995 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2996 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2997 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2998 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2999 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |

|     | ... | 37668 | 37669 | 37670 | 37671 | 37672 | 37673 | 37674 | 37675 | 37676 | \ |
|-----|-----|-------|-------|-------|-------|-------|-------|-------|-------|-------|---|

```
0      ...    0    0    0    0    0    0    0    0    0
1      ...    0    0    0    0    0    0    0    0    0
2      ...    0    0    0    0    0    0    0    0    0
3      ...    0    0    0    0    0    0    0    0    0
4      ...    0    0    0    0    0    0    0    0    0

...    ...    ...  ...  ...  ...  ...  ...  ...  ...  ...
2995   ...    0    0    0    0    0    0    0    0    0
2996   ...    0    0    0    0    0    0    0    0    0
2997   ...    0    0    0    0    0    0    0    0    0
2998   ...    0    0    0    0    0    0    0    0    0
2999   ...    0    0    0    0    0    0    0    0    0

        37677
0            0
1            0
2            0
3            0
4            0

...        ...
2995         0
2996         0
2997         0
2998         0
2999         0

[3000  rows x  37678 columns]
```

[46] : `vect_data.dtypes`

```
[46] : 0          int64
       1          int64
       2          int64
       3          int64
       4          int64

                  ...
       37673      int64
       37674      int64
       37675      int64
       37676      int64
       37677      int64
       Length: 37678, dtype: object
```

[47] :
```
x1 = vect_data
y = samp['True']
x_train,x_test,y_train,y_test  =  train_test_split(x1,y,  test_size  =  0.2)
```

```
[48] :  ES = tf.keras.callbacks.EarlyStopping(
            min_delta = 0.001 ,
            patience = 15 ,
            restore_best_weights = True
        )
        model = tf.keras.Sequential([
            layers.BatchNormalization(),
            layers.Dropout(0.3),
            layers.Dense(100, activation = 'relu'),
            layers.Dense(1, activation='sigmoid')
        ])
        model.compile(
            optimizer = 'adam',
            metrics = ['binary_accuracy'],
            loss = 'binary_crossentropy'
        )
        truth = model.fit(tf.cast(x_train , tf.float32),y_train,
                          validation_data =(x_test,y_test),
                          verbose = 0,
                        callbacks = [ES],
                          batch_size = 100,
                        epochs = 500)
        history_df = pd.DataFrame(truth.history)
        history_df.loc[:, ['loss', 'val_loss']].plot();
```