



4222-SURYA GROUP OF INSTITUTION



VIKIRAVANDI-605 652.

NAAN MUDHALVAN PROJECT

CREATE CHATBOT IN PYTHON

PHASE 3: DEVELOPMENT PART 1

Presented by:

S.SANTHOSH

REG NO:422221106308

ECE DEPARTMENT

AI_PHASE 3:

Some common steps in data preprocessing include:

Data preprocessing is an important step in the data mining process that involves cleaning and transforming raw data to make it suitable for analysis. Some common steps in data preprocessing include

Data Cleaning:

This involves identifying and correcting errors or inconsistencies in the data, such as missing values, outliers, and duplicates. Various techniques can be used for data cleaning, such as imputation, removal, and transformation.

Data Integration:

This involves combining data from multiple sources to create a unified dataset. Data integration can be challenging as it requires handling data with different formats, structures, and semantics. Techniques such as record linkage and data fusion can be used for data integration.

Data Transformation:

This involves converting the data into a suitable format for analysis. Common techniques used in data transformation include normalization, standardization, and discretization. Normalization is used to scale the data to a common range, while standardization is used to transform the data to have zero mean and unit variance. Discretization is used to convert continuous data into discrete categories.

Data Reduction:

This involves reducing the size of the dataset while preserving the important information. Data reduction can be achieved through techniques such as feature selection and feature extraction. Feature selection involves selecting a subset of relevant features from the dataset, while feature extraction involves transforming the data into a lower-dimensional space while preserving the important information.

Data Discretization:

This involves dividing continuous data into discrete categories or intervals. Discretization is often used in data mining and machine learning algorithms that require categorical data. Discretization can be achieved through techniques such as equal width binning, equal frequency binning, and clustering.

Data Normalization:

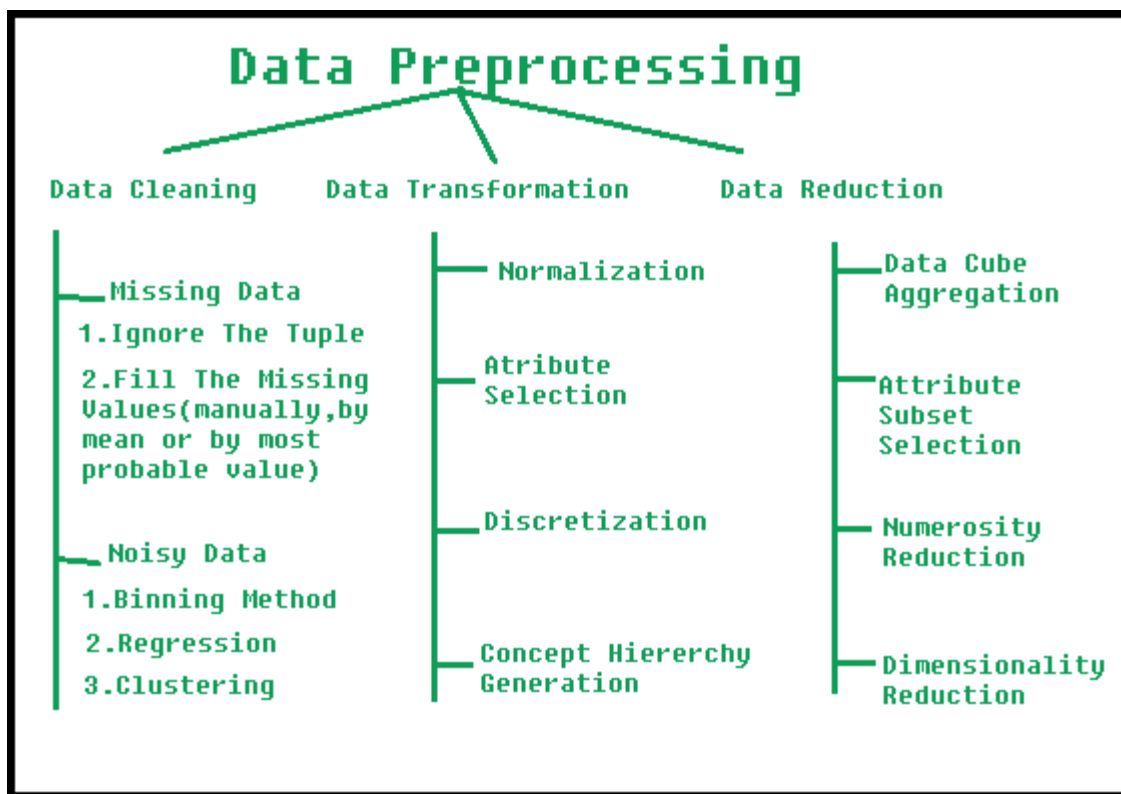
This involves scaling the data to a common range, such as between 0 and 1 or -1 and 1. Normalization is often used to handle data with different units and scales. Common normalization techniques include min-max normalization, z-score normalization, and decimal scaling.

Data preprocessing plays a crucial role in ensuring the quality of data and the accuracy of the analysis results. The specific steps involved in data preprocessing may vary depending on the nature of the data and the analysis goals.

By performing these steps, the data mining process becomes more efficient and the results become more accurate

Preprocessing in Data Mining:

Data preprocessing is a data mining technique which is used to transform the raw data in a useful and efficient format.



Steps Involved in Data Preprocessing:

Data Cleaning:

The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data, noisy data etc.

Missing Data:

This situation arises when some data is missing in the data. It can be handled in various ways.

Some of them are:

Ignore the tuples:

This approach is suitable only when the dataset we have is quite large and multiple values are missing within a tuple.

Fill the Missing values:

There are various ways to do this task. You can choose to fill the missing values manually, by attribute mean or the most probable value.

Noisy Data:

Noisy data is a meaningless data that can't be interpreted by machines. It can be generated due to faulty data collection, data entry errors etc. It can be handled in following ways

1.Binning Method:

This method works on sorted data in order to smooth it. The whole data is divided into segments of equal size and then various methods are performed to complete the task. Each segment is handled separately. One can replace all data in a segment by its mean or boundary values can be used to complete the task.

2.Regression:

Here data can be made smooth by fitting it to a regression function. The regression used may be linear (having one independent variable) or multiple (having multiple independent variables)

3.Clustering:

This approach groups the similar data in a cluster. The outliers may be undetected or it will fall outside the clusters.

Data Transformation:

This step is taken in order to transform the data in appropriate forms suitable for mining process. This involves following ways

Normalization:

It is done in order to scale the data values in a specified range (1.0 to 1.0 or 0.0 to 1.0)

Attribute Selection:

In this strategy, new attributes are constructed from the given set of attributes to help the mining process.

Discretization:

This is done to replace the raw values of numeric attribute by interval levels or conceptual levels.

Concept Hierarchy Generation:

Here attributes are converted from lower level to higher level in hierarchy. For Example-The attribute “city” can be converted to “country”.

Data Preprocessing:

Data Visualization

```
df['question tokens']=df['question'].apply(lambda x:len(x.split()))
```

```
df['answer tokens']=df['answer'].apply(lambda x:len(x.split()))
```

```
plt.style.use('fivethirtyeight')
```

```
fig,ax=plt.subplots(nrows=1,ncols=2,figsize=(20,5))
```

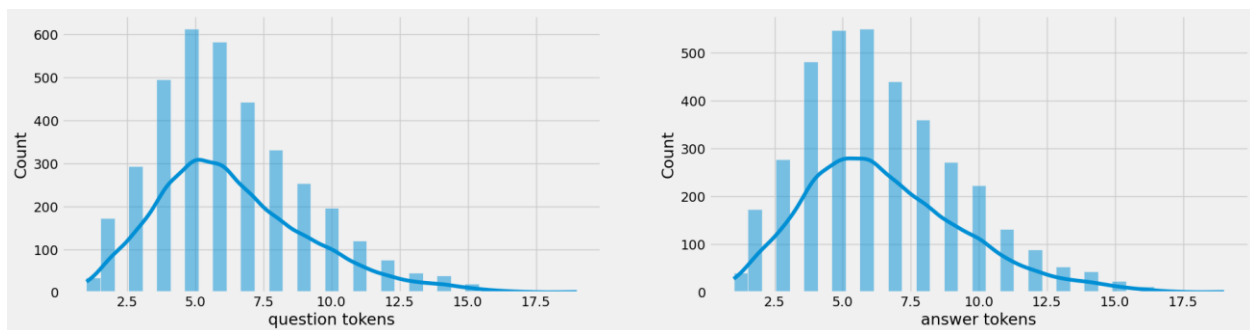
```
sns.set_palette('Set2')
```

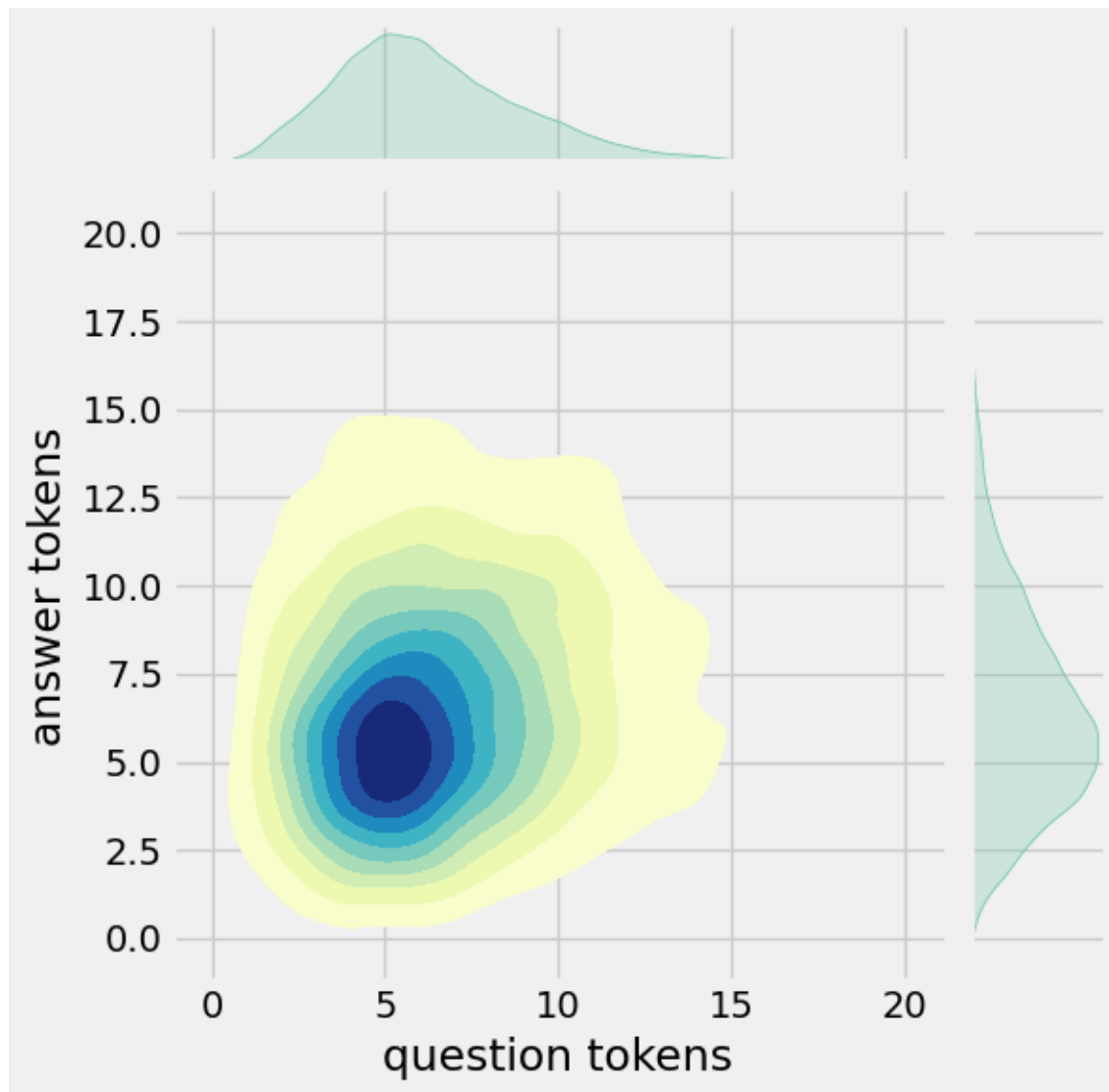
```
sns.histplot(x=df['question tokens'],data=df,kde=True,ax=ax[0])
```

```
sns.histplot(x=df['answer tokens'],data=df,kde=True,ax=ax[1])
```

```
sns.jointplot(x='question tokens',y='answer  
tokens',data=df,kind='kde',fill=True,cmap='YlGnBu')
```

```
plt.show()
```





Text Cleaning

```
def clean_text(text):  
    text=re.sub('-', ' ',text.lower())  
    text=re.sub('[.]', ' . ',text)  
    text=re.sub('[1]', ' 1 ',text)  
    text=re.sub('[2]', ' 2 ',text)  
    text=re.sub('[3]', ' 3 ',text)
```



```
text=re.sub('[4]',' 4 ',text)
text=re.sub('[5]',' 5 ',text)
text=re.sub('[6]',' 6 ',text)
text=re.sub('[7]',' 7 ',text)
text=re.sub('[8]',' 8 ',text)
text=re.sub('[9]',' 9 ',text)
text=re.sub('[0]',' 0 ',text)
text=re.sub('[,]',' , ',text)
text=re.sub('[?]' ,' ? ',text)
text=re.sub('[!]' ,' ! ',text)
text=re.sub('[\$]' ,' $ ',text)
text=re.sub('[&]' ,' & ',text)
text=re.sub('[/]' ,' / ',text)
text=re.sub('[:]' ,' : ',text)
text=re.sub('[:,]' ,' ; ',text)
text=re.sub('[*]' ,' * ',text)
text=re.sub('[\\]' ,' \' ',text)
text=re.sub('[\"]' ,' \" ',text)
text=re.sub('\\t',' ',text)
return text
```

```
df.drop(columns=['answer tokens','question tokens'],axis=1,inplace=True)
```

```
df['encoder_inputs']=df['question'].apply(clean_text)
df['decoder_targets']=df['answer'].apply(clean_text)+' <end>'
df['decoder_inputs']='<start> '+df['answer'].apply(clean_text)+' <end>'
df.head(10)
```

After preprocessing: for example , if your birth date is january 1 2 , 1 9 8 7 , write 0 1 / 1 2 / 8 7 .

Max encoder input length: 27

Max decoder input length: 29

Max decoder target length: 28

Tokenization

```
vectorize_layer=TextVectorization(
    max_tokens=vocab_size,
    standardize=None,
    output_mode='int',
    output_sequence_length=max_sequence_length
)
vectorize_layer.adapt(df['encoder_inputs']+' '+df['decoder_targets']+'
<start> <end>')
vocab_size=len(vectorize_layer.get_vocabulary())
print(f'Vocab size: {len(vectorize_layer.get_vocabulary())}')
print(f'{vectorize_layer.get_vocabulary()[:12]}')
```

Vocab size: 2443

```
[' ', '[UNK]', '<end>', '!', '<start>', '"', 'i', '?', 'you', ',', 'the', 'to']
```

```
def sequences2ids(sequence):
```

```
    return vectorize_layer(sequence)
```

```
def ids2sequences(ids):
```

```
    decode=""
```

```
    if type(ids)==int:
```

```
        ids=[ids]
```

```
    for id in ids:
```

```
        decode+=vectorize_layer.get_vocabulary()[id]+' '
```

```
    return decode
```

```
x=sequences2ids(df['encoder_inputs'])
```

```
yd=sequences2ids(df['decoder_inputs'])
```

```
y=sequences2ids(df['decoder_targets'])
```

```
print(f'Question sentence: hi , how are you ?')
```

```
print(f'Question to tokens: {sequences2ids("hi , how are you ?")[:10]}')
```

```
print(f'Encoder input shape: {x.shape}')
```

```
print(f'Decoder input shape: {yd.shape}')
```

```
print(f'Decoder target shape: {y.shape}')
```

```
Question sentence: hi , how are you ?
```

Question to tokens: [1971 9 45 24 8 7 0 0 0 0]

Encoder input shape: (3725, 30)

Decoder input shape: (3725, 30)

Decoder target shape: (3725, 30)

```
print(f'Encoder input: {x[0][:12]} ...')
```

```
print(f'Decoder input: {yd[0][:12]} ...') # shifted by one time step of  
the target as input to decoder is the output of the previous timestep
```

```
print(f'Decoder target: {y[0][:12]} ...')
```

```
Encoder input: [1971 9 45 24 8 194 7 0 0 0 0 0] ...
```

```
Decoder input: [ 4 6 5 38 646 3 45 41 563 7 2 0] ...
```

```
Decoder target: [ 6 5 38 646 3 45 41 563 7 2 0 0] ...
```

```
data=tf.data.Dataset.from_tensor_slices((x,yd,y))
```

```
data=data.shuffle(buffer_size)
```

```
train_data=data.take(int(.9*len(data)))
```

```
train_data=train_data.cache()
```

```
train_data=train_data.shuffle(buffer_size)
```

```
train_data=train_data.batch(batch_size)
```

```
train_data=train_data.prefetch(tf.data.AUTOTUNE)
```

```
train_data_iterator=train_data.as_numpy_iterator()
```

```
val_data=data.skip(int(.9*len(data))).take(int(.1*len(data)))
```

```
val_data=val_data.batch(batch_size)
```

```
val_data=val_data.prefetch(tf.data.AUTOTUNE)

_=train_data_iterator.next()

print(f'Number of train batches: {len(train_data)}')
print(f'Number of training data: {len(train_data)*batch_size}')
print(f'Number of validation batches: {len(val_data)}')
print(f'Number of validation data: {len(val_data)*batch_size}')
print(f'Encoder Input shape (with batches): {_[0].shape}')
print(f'Decoder Input shape (with batches): {_[1].shape}')
print(f'Target Output shape (with batches): {_[2].shape}')

Number of train batches: 23
Number of training data: 3427
Number of validation batches: 3
Number of validation data: 447
Encoder Input shape (with batches): (149, 30)
Decoder Input shape (with batches): (149, 30)
Target Output shape (with batches): (149, 30)
```

Data Reduction:

Data reduction is a crucial step in the data mining process that involves reducing the size of the dataset while preserving the important information. This is done to improve the efficiency of data analysis and to avoid overfitting of the model. Some common steps involved in data reduction are:

Feature Selection

This involves selecting a subset of relevant features from the dataset. Feature selection is often performed to remove irrelevant or redundant features from the dataset. It can be done using various techniques such as correlation analysis, mutual information, and principal component analysis (PCA).

Sampling

This involves selecting a subset of data points from the dataset. Sampling is often used to reduce the size of the dataset while preserving the important information. It can be done using techniques such as random sampling, stratified sampling, and systematic sampling.

Clustering

This involves grouping similar data points together into clusters. Clustering is often used to reduce the size of the dataset by replacing similar data points with a representative centroid. It can be done using techniques such as k-means, hierarchical clustering, and density-based clustering.

CONCLUSION

A chatbot is one of the simple ways to transport data from a computer without having to think for proper keywords to look up in a search or browse several web pages to collect information; users can easily type their query in natural language and retrieve information. In this paper, information about the design, implementation of the chatbot has been presented. From the survey above, it can be said that the development and improvement of chatbot design grow at an unpredictable rate due to variety of methods and approaches used to design a chatbot. Chatbot is a great tool for quick interaction with the user.