# MCQ GENERATION

## Using Machine Learning

# TABLE OF CONTENTS

# DISCLAIMER

All software and hardware used or referenced in this guide belong to their respective vendor. We developed this guide based on our development infrastructure and this guide may or may not work on others systems and technical infrastructure. We are not liable for any direct or indirect problems caused by users using this guide.

# EXECUTIVE SUMMARY

The purpose of this document is to provide adequate information to users to implement a Supervised Machine Learning model. In order to achieve this, we are using one of the most common problem that occurs at Educational Institutions. Traditionally, a Teacher picks a few questions from the Chapter and these questions are repeated every year. Each Chapter will require its own way of Learning and Understanding. In order to generate questions in all Dimensions Possible for every Chapter, we use Machine Learning Techniques, so that Machine can Learn, Understand and Frame as many questions as possible instead of a Teacher taking questions by themselves.

# BUSINESS PROBLEM

# PROBLEM STATEMENT

Given the Paragraph, Frame Multiple Choice Questions and Generate Distractors (Options) similar to the Answer.

# BUSINESS CHALLENGES

- Takes more time
- Human intensive
- The questions generated by humans are limited to predefined context

# BUSINESS CONTEXT

In an Educational Institution, a teacher picks few Questions from a Chapter. These Questions are repeated every Semester, for example, a Maths Teacher creates few types of Questions for a chapter, it is repeated for the students passing that semester for many years. The learning speed of each student varies, there can't be common questions across all the students. So the questions has to be dynamically driven. In order to achieve this, we are using Machine Learning Technique, so that the Machines can learn, understand and create Questions in all Dimensions Possible by itself instead of a Teacher picking questions by themselves.

# DATA MANAGEMENT

**There are three types of data sets Training, Test and Dev that are used at various stage of Implementation. Training dataset is the largest of three of them, while test data functions as seal of approval and you don't need to use till the end of the development.**
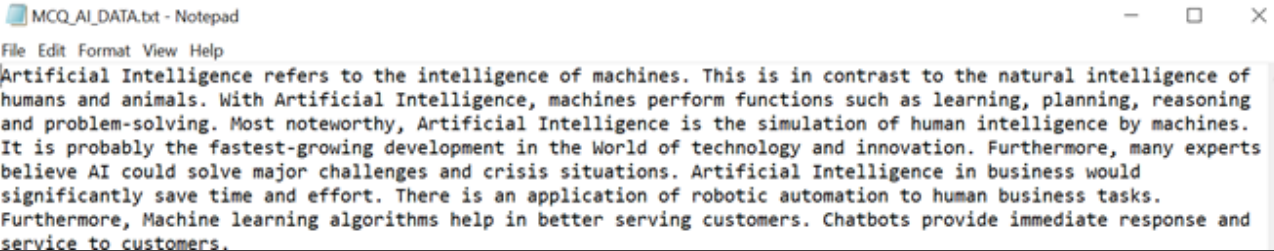
# TRAINING DATASET

The training data set is the actual dataset used to train the model for performing various Machine Learning Operations (Regression, Classification, Clustering etc.). This is the actual data with which the models learn with various API and algorithm to train the machine to work automatically.

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | title | context | question | answer | answer_st | answer_e | answer_st | answer_end_word_idx | | |
| 2 | University | Architectu | To whom | Saint Bern | 515 | 541 | 102 | 104 | | |
| 3 | University | Architectu | What is in | a copper s | 188 | 213 | 37 | 41 | | |
| 4 | University | Architectu | The Basilic | the Main I | 279 | 296 | 57 | 59 | | |
| 5 | University | Architectu | What is th | a Marian | 381 | 420 | 76 | 82 | | |
| 6 | University | Architectu | What sits | a golden s | 92 | 126 | 17 | 23 | | |
| 7 | University | As at most | When did | Septembe | 248 | 262 | 45 | 46 | | |
| 8 | University | As at most | How ofter | twice | 441 | 446 | 78 | 78 | | |
| 9 | University | As at most | What is th | The Obser | 598 | 610 | 216 | 217 | | |
| 10 | University | As at most | How many | three | 126 | 131 | 222 | 222 | | |
| 11 | University | As at most | In what ye | 1987 | 908 | 912 | 158 | 158 | | |
| 12 | University | The unive | Where is t | Rome | 119 | 123 | 22 | 22 | | |
| 13 | University | The unive | What is th | Moreau S | 145 | 160 | 121 | 122 | | |
| 14 | University | The unive | What is th | Old Colleg | 234 | 245 | 46 | 47 | | |
| 15 | University | The unive | What indi | Retired pr | 356 | 384 | 68 | 71 | | |
| 16 | University | The unive | Which pri: | Buechner | 675 | 703 | 125 | 128 | | |
| 17 | University | The Colleg | How many | eight | 487 | 492 | 79 | 79 | | |
| 18 | University | The Colleg | In what ye | 1920 | 46 | 50 | 7 | 7 | | |
| 19 | University | The Colleg | Before the | the Colleg | 126 | 148 | | | | |

```
the -0.038194 -0.24487 0.72812 -0.39961 0.083172 0.043953 -0.39141 0.3344 -0.57545 0.087459 0.
, -0.10767 0.11053 0.59812 -0.54361 0.67396 0.10663 0.038867 0.35481 0.06351 -0.094189 0.15786
. -0.33979 0.20941 0.46348 -0.64792 -0.38377 0.038034 0.17127 0.15978 0.46619 -0.019169 0.41471
of -0.1529 -0.24279 0.89837 0.16996 0.53516 0.48784 -0.58826 -0.17982 -1.3581 0.42541 0.15377 (
to -0.1897 0.050024 0.19084 -0.049184 -0.089737 0.21006 -0.54952 0.098377 -0.20135 0.34241 -0.(
and -0.071953 0.23127 0.023731 -0.50638 0.33923 0.1959 -0.32943 0.18364 -0.18057 0.28963 0.204
in 0.085703 -0.22201 0.16569 0.13373 0.38239 0.35401 0.01287 0.22461 -0.43817 0.50164 -0.35874
a -0.27086 0.044006 -0.02026 -0.17395 0.6444 0.71213 0.3551 0.47138 -0.29637 0.54427 -0.72294
" -0.30457 -0.23645 0.17576 -0.72854 -0.28343 -0.2564 0.26587 0.025309 -0.074775 -0.3766 -0.05
's 0.58854 -0.2025 0.73479 -0.68338 -0.19675 -0.1802 -0.39177 0.34172 -0.60561 0.63816 -0.2669
for -0.14401 0.32554 0.14257 -0.099227 0.72536 0.19321 -0.24188 0.20223 -0.89599 0.15215 0.035
- -1.2557 0.61036 0.56793 -0.96596 -0.45249 -0.071696 0.57122 -0.31292 -0.43814 0.90622 0.0696
that -0.093337 0.19043 0.68457 -0.41548 -0.22777 -0.11803 -0.095434 0.19613 0.17785 -0.020244
on -0.21863 -0.42664 0.5196 0.0043103 0.58045 -0.10873 -0.37726 0.4566 -0.60627 -0.075773 0.11
is -0.54264 0.41476 1.0322 -0.40244 0.46691 0.21816 -0.074864 0.47332 0.080996 -0.22079 -0.128
was 0.13717 -0.54287 0.19419 -0.29953 0.17545 0.084672 0.67752 0.098295 -0.035611 0.21334 0.51
said -0.13128 -0.452 0.043399 -0.99798 -0.21053 -0.95868 -0.24609 0.48413 0.18178 0.475 -0.223
with -0.43608 0.39104 0.51657 -0.13861 0.2029 0.50723 -0.012544 0.22948 -0.6316 0.21199 -0.018
he 0.1225 -0.058833 0.23658 -0.28877 -0.028181 0.31524 0.070229 0.16447 -0.027623 0.25214 0.21
```

# TEST DATASET

Test data set helps you to validate that the training has happened efficiently in terms of either accuracy, or precision so on. Actually, such data is used for testing the model whether it is responding or working appropriately or not.

MCQ_AI_DATA.txt - Notepad

File Edit Format View Help

Artificial Intelligence refers to the intelligence of machines. This is in contrast to the natural intelligence of humans and animals. With Artificial Intelligence, machines perform functions such as learning, planning, reasoning and problem-solving. Most noteworthy, Artificial Intelligence is the simulation of human intelligence by machines. It is probably the fastest-growing development in the World of technology and innovation. Furthermore, many experts believe AI could solve major challenges and crisis situations. Artificial Intelligence in business would significantly save time and effort. There is an application of robotic automation to human business tasks. Furthermore, Machine learning algorithms help in better serving customers. Chatbots provide immediate response and service to customers.

DeepSphere.AI
Enterprise AI and IIoT for Analytics

# MACHINE LEARNING LIBRARIES USED

| 1 | spaCy |
|---|-------|

| 2 | Pandas |
|---|--------|

| 3 | Gensim |
|---|--------|

| 4 | Random |
|---|--------|

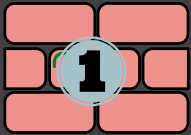| 5 | Pickle |
|---|--------|

# MODEL USED:

## Gaussian Naive Bayes

# MODEL BUILDING BLOCKS

There are several technical and functional components involved in implementing this model. Here are the key building blocks to implement the model.

# MACHINE LEARNING BUILDING BLOCKS

1 GOOGLE CLOUD

PYTHON 2

3 JUPYTER NOTEBOOK

PANDAS 4

5 PYSPARK - DATA ENGINEERING

GENSIM 6

7 PICKLE

RANDOM 8

9 SPACY

DeepSphere.AI
Enterprise AI and IIoT for Analytics

# MACHINE LEARNING IMPLEMENTATION STEPS

Here are the key steps that are involved to implement a deep learning model. You can customize these steps as needed and we have developed these steps for learning-purpose only.

# Implementation Steps

**Step 1**
Setting up Development Infrastructure

**Step 2**
Data Preprocessing

**Step 3**
Feature Engineering

**Step 4**
Data Preparation

**Step 5**
Finalising the ML Components

**Step 6**
Train ML

**Step 7**
Generating Distractors

**Step 8**
Test ML

**Step 9**
Verify the Model Outcome

**Step 10**
Retrain , Retest and Redeploy

DeepSphere.AI
Enterprise AI and IIoT for Analytics

# ML MODEL IMPLEMENTATION STEPS

# STEP 1-SETTING UP DEVELOPMENT INFRASTRUCTURE:

For our Model Implementation we need theFollowing Libraries:

**Pandas:**Pandas is a library used for data manipulation and analysis. For Our Implementation we are using it for Importing the Data file & Creating Dataframes (Stores the Data).

**Spacy**: It can be used to build information extraction or natural language understanding systems, or to pre-process text for deep learning. It can be used to build information extraction or natural language understanding systems, or to pre-process text for deep learning.

**Gensim**: Gensim is a Python library for topic modelling, document indexing and similarity retrieval with large corpora.

**Random**: Python offers random module that can generate random numbers. These are pseudo-random number
as the sequence of number generated depends on the seed. If the seeding value is same, the sequence will be the same.

```
                mcq.py                    •
 1   /**********************************************************************
 2
 3   FILE NAME      : mcq_question_generation.py
 4   Purpose        : To generate a paragraph of text automatically
 5   Author         : DeepSphere.AI, Inc.
 6   Date and Time  : 12/07/2020 10:00 hrs
 7   Version        : 1.0
 8
 9   /**********************************************************************
10
11   ###############################################
12   # Step 1 - Import the Required Libraries     #
13   ###############################################
14
15   import pandas as pd
16   from IPython.display import Markdown, display, clear_output
17   import spacy
18   from spacy import displacy
19   import _pickle as cPickle
20   from pathlib import Path
21   import gensim
22   from gensim.test.utils import datapath, get_tmpfile
23   from gensim.models import KeyedVectors
24   import random
25
26   ###############################################
27   #    Step 2-Import Clustering Data           #
28   ###############################################
29   vAR_Pickle_Data1 =C:\AI\AUTOMATIC QUESTION GENERATION\ML\TRAINING DATA\pickles\nb-predictor.pkl
30   vAR_Pickle_Data2 =C:\AI\AUTOMATIC QUESTION GENERATION\ML\TRAINING DATA\pickles\wordsDf.pkl
31   vAR_Training_Data1 = C:\AI\AUTOMATIC QUESTION GENERATION\ML\TRAINING DATA\embeddings\glove.6B.300d.txt
32   vAR_Training_Data2 =C:\AI\AUTOMATIC QUESTION GENERATION\ML\TRAINING DATA\embeddings\word2vec-glove.6B.300d.txt
33   vAR_Test_Data = C:\AI\AUTOMATIC QUESTION GENERATION\ML\TEST DATA\MCQ_AI_DATA.txt
34
```

**2**

# STEP 2 - DATA PREPROCESSING

Next immediate step after importing all
libraries is Data preprocessing i.e. Pickling. "Pickling"
is the process whereby a Python object hierarchy is converted
into a byte stream.

```python
35    ###############################################
36    #              Step 3-Pickling                #
37    ###############################################
38
39    def dumpPickle(fileName, content):
40        pickleFile = open(fileName, 'wb')
41        cPickle.dump(content, pickleFile, -1)
42        pickleFile.close()
43
44    def loadPickle(fileName):
45        file = open(fileName, 'rb')
46        content = cPickle.load(file)
47        file.close()
48        return content
49
50    def pickleExists(fileName):
51        file = Path(fileName)
52
53        if file.is_file():
54            return True
55
56        return False
57
```

DeepSphere.AI
Enterprise AI and IIoT for Analytics

3

# STEP 3 - FEATURE ENGINEERING

Step 3 of the Implementation is Feature Generation/ Feature Engineering. Machine learning works on a simple rule – if you put garbage in, you will only get garbage to come out. By garbage here, we mean noise in data. This becomes even more important when the numbers of features are very large. We need only those features (Input) that are function of the Labels (Outputs). Ex: To Predict whether the given fruit is an apple or orange Color/Texture of the Fruit becomes a feature to be Considered. If the Color/Texture is Red then it an Apple, If it's Orange its Orange.

```python
58   ########################################################################
59   #              Step 4-Extract Words and Generate Features              #
60   ########################################################################
61
62   vAR_model = LinearRegression()
63   vAR_model.fit(vAR_Features_Train,vAR_Labels_Train)
64   import en_core_web_sm
65   nlp = spacy.load('en_core_web_sm')
66
67   ###############################
68   #Extract answers and the sentence they are in
69   ###############################
70   def extractAnswers(qas, doc):
71       answers = []
72
73       senStart = 0
74       senId = 0
75
76       for sentence in doc.sents:
77           senLen = len(sentence.text)
78
79           for answer in qas:
80               answerStart = answer['answers'][0]['answer_start']
81
82               if (answerStart >= senStart and answerStart < (senStart + senLen)):
83                   answers.append({'sentenceId': senId, 'text': answer['answers'][0]['text']})
84
85           senStart += senLen
86           senId += 1
87
88       return answers
89   ###############################
90   # Cleaning answers from stopwords
91   ###############################
92   def tokenIsAnswer(token, sentenceId, answers):
93       for i in range(len(answers)):
94           if (answers[i]['sentenceId'] == sentenceId):
95               if (answers[i]['text'] == token):
96                   return True
97       return False
98
99
```

DeepSphere.AI
Enterprise AI and IIoT for Analytics

# 4

# STEP 4 - DATA PREPERATION

Step 4 involved fixing named entities start points. In information extraction, a named entity is a real-world object, such as persons, locations, organizations, products, etc., that can be denoted with a proper name. It can be abstract or have a physical existence. Hence, These named entities can be used to select potential Questions and answers.

```python
100    ################################################
101    #     Step 5-Fixing named entities start points #
102    ################################################
103
104    #################################
105    #Save named entities start points
106    #################################
107
108    def getNEStartIndexs(doc):
109        neStarts = {}
110        for ne in doc.ents:
111            neStarts[ne.start] = ne
112
113        return neStarts
114
115    def getSentenceStartIndexes(doc):
116        senStarts = []
117
118        for sentence in doc.sents:
119            senStarts.append(sentence[0].i)
120
121        return senStarts
122
123    def getSentenceForWordPosition(wordPos, senStarts):
124        for i in range(1, len(senStarts)):
125            if (wordPos < senStarts[i]):
126                return i - 1
127
```

```python
128    def addWordsForParagrapgh(newWords, text):
129        doc = nlp(text)
130
131        neStarts = getNEStartIndexs(doc)
132        senStarts = getSentenceStartIndexes(doc)
133
134        #index of word in spacy doc text
135        i = 0
136
137        while (i < len(doc)):
138            #If the token is a start of a Named Entity, add it and push to index to end of the NE
139            if (i in neStarts):
140                word = neStarts[i]
141                #add word
142                currentSentence = getSentenceForWordPosition(word.start, senStarts)
143                wordLen = word.end - word.start
144                shape = ''
145                for wordIndex in range(word.start, word.end):
146                    shape += (' ' + doc[wordIndex].shape_)
147
148                newWords.append([word.text,
149                                 0,
150                                 0,
151                                 currentSentence,
152                                 wordLen,
153                                 word.label_,
154                                 None,
155                                 None,
156                                 None,
157                                 shape])
158            i = neStarts[i].end - 1
159            #If not a NE, add the word if it's not a stopword or a non-alpha (not regular letters)
160            else:
161                if (doc[i].is_stop == False and doc[i].is_alpha == True):
162                    word = doc[i]
163
164                    currentSentence = getSentenceForWordPosition(i, senStarts)
165                    wordLen = 1
166
167                    newWords.append([word.text,
168                                     0,
169                                     0,
170                                     currentSentence,
```

```
157                             shape])
158              i = neStarts[i].end - 1
159          #If not a NE, add the word if it's not a stopword or a non-alpha (not regular letters)
160          else:
161              if (doc[i].is_stop == False and doc[i].is_alpha == True):
162                  word = doc[i]
163
164                  currentSentence = getSentenceForWordPosition(i, senStarts)
165                  wordLen = 1
166
167                  newWords.append([word.text,
168                                  0,
169                                  0,
170                                  currentSentence,
171                                  wordLen,
172                                  None,
173                                  word.pos_,
174                                  word.tag_,
175                                  word.dep_,
176                                  word.shape_])
177          i += 1
178
179  def oneHotEncodeColumns(df):
180      columnsToEncode = ['NER', 'POS', "TAG", 'DEP']
181
182      for column in columnsToEncode:
183          one_hot = pd.get_dummies(df[column])
184          one_hot = one_hot.add_prefix(column + '_')
185
186          df = df.drop(column, axis = 1)
187          df = df.join(one_hot)
188
189      return df
190
```

# STEP 5 - ANALYSING AND FINALISING THE ML COMPONENTS

As a next step we need to predict whether a word is a keyword. Here we do one-hot Encoding, Drop Unused Columns and add the missing Columns.

```
191    #################################################
192    #     Step 6-Predict whether word is a keyword  #
193    #################################################
194
195    def generateDf(text):
196        words = []
197        addWordsForParagrapgh(words, text)
198
199        wordColums = ['text', 'titleId', 'paragrapghId', 'sentenceId','wordCount', 'NER', 'POS', 'TAG', 'DEP','shape']
200        df = pd.DataFrame(words, columns=wordColums)
201
202        return df
203    def prepareDf(df):
204        #One-hot encoding
205        wordsDf = oneHotEncodeColumns(df)
206
207        #Drop unused columns
208        columnsToDrop = ['text', 'titleId', 'paragrapghId', 'sentenceId', 'shape']
209        wordsDf = wordsDf.drop(columnsToDrop, axis = 1)
210
211        #Add missing colums
212        predictorColumns = ['wordCount','NER_CARDINAL','NER_DATE','NER_EVENT','NER_FAC','NER_GPE','NER_LANGUAGE','NER_LAW','NER_LOC','NER_MONEY','NER_NORP
     ','TAG_.','TAG_ADD','TAG_AFX','TAG_CC','TAG_CD','TAG_DT','TAG_EX','TAG_FW','TAG_IN','TAG_JJ','TAG_JJR','TAG_JJS','TAG_LS','TAG_MD','TAG_NFP','TAG_
     x','DEP_auxpass','DEP_case','DEP_cc','DEP_ccomp','DEP_compound','DEP_conj','DEP_csubj','DEP_csubjpass','DEP_dative','DEP_dep','DEP_det','DEP_dobj'
213
214        for feature in predictorColumns:
215            if feature not in wordsDf.columns:
216                wordsDf[feature] = 0
217
218        return wordsDf
219    def predictWords(wordsDf, df):
220
221        predictorPickleName = vAR_Pickle_Data1
222        predictor = loadPickle(predictorPickleName)
223
224        y_pred = predictor.predict_proba(wordsDf)
225
226        labeledAnswers = []
227        for i in range(len(y_pred)):
228            labeledAnswers.append({'word': df.iloc[i]['text'], 'prob': y_pred[i][0]})
229
230        return labeledAnswers
231
```

# 6

# STEP 6 - TRAIN ML

By step 6, We Extract Questions from the Potential Sentence. We also group questions and answers

```
233    ################################################
234    #            Step 7-Extract Questions          #
235    ################################################
236
237    vAR_Labels_Pred = vAR_model.predict(vAR_Features_Test).astype(int)
238    def blankAnswer(firstTokenIndex, lastTokenIndex, sentStart, sentEnd, doc):
239        leftPartStart = doc[sentStart].idx
240        leftPartEnd = doc[firstTokenIndex].idx
241        rightPartStart = doc[lastTokenIndex].idx + len(doc[lastTokenIndex])
242        rightPartEnd = doc[sentEnd - 1].idx + len(doc[sentEnd - 1])
243
244        question = doc.text[leftPartStart:leftPartEnd] + '_____' + doc.text[rightPartStart:rightPartEnd]
245
246        return question
```

```
248  ###############################################################
249  #            Step 8-Grouping Questions and Answers            #
250  ###############################################################
251
252  def addQuestions(answers, text):
253      doc = nlp(text)
254      currAnswerIndex = 0
255      qaPair = []
256
257      #Check wheter each token is the next answer
258      for sent in doc.sents:
259          for token in sent:
260
261              #If all the answers have been found, stop looking
262              if currAnswerIndex >= len(answers):
263                  break
264
265              #In the case where the answer is consisted of more than one token, check the following tokens as well.
266              answerDoc = nlp(answers[currAnswerIndex]['word'])
267              answerIsFound = True
268
269              for j in range(len(answerDoc)):
270                  if token.i + j >= len(doc) or doc[token.i + j].text != answerDoc[j].text:
271                      answerIsFound = False
272
273              #If the current token is corresponding with the answer, add it
274              if answerIsFound:
275                  question = blankAnswer(token.i, token.i + len(answerDoc) - 1, sent.start, sent.end, doc)
276
277                  qaPair.append({'question' : question, 'answer': answers[currAnswerIndex]['word'], 'prob': answers[currAnswerIndex]['prob']})
278
279                  currAnswerIndex += 1
280
281      return qaPair
282  def sortAnswers(qaPairs):
283      orderedQaPairs = sorted(qaPairs, key=lambda qaPair: qaPair['prob'])
284
285      return orderedQaPairs
286
```

**8**

# STEP 8 - GENERATING DISTRACTORS

Next, we do a very important step to test the knowledge. We generate very similar options to answer. This is a very crucial step in creating MCQ Questions.

```
287  ########################################################
288  #            Step 9-Generationg Distractors           #
289  ########################################################
290
291  glove_file = vAR_Training_Data1
292  tmp_file = vAR_Training_Data2
293
294  from gensim.scripts.glove2word2vec import glove2word2vec
295  glove2word2vec(glove_file, tmp_file)
296  model = KeyedVectors.load_word2vec_format(tmp_file)
297  def generate_distractors(answer, count):
298      answer = str.lower(answer)
299
300      ##Extracting closest words for the answer.
301      try:
302          closestWords = model.most_similar(positive=[answer], topn=count)
303      except:
304          #In case the word is not in the vocabulary, or other problem not loading embeddings
305          return []
306
307      #Return count many distractors
308      distractors = list(map(lambda x: x[0], closestWords))[0:count]
309
310      return distractors
311  def addDistractors(qaPairs, count):
312      for qaPair in qaPairs:
313          distractors = generate_distractors(qaPair['answer'], count)
314          qaPair['distractors'] = distractors
315
316      return qaPairs
317
```

# 9

## STEP 9 - TEST ML

Yes, This is our MAIN function which actually generates Questions. Here we integrate all the functions to end our Model.

```
318    ###############################################
319    #              Step 10-Main Function          #
320    ###############################################
321
322    def generateQuestions(text, count):
323
324        ################################
325        # Extract words
326        ################################
327        df = generateDf(text)
328        wordsDf = prepareDf(df)
329
330        ################################
331        # Predict
332        ################################
333        labeledAnswers = predictWords(wordsDf, df)
334
335        ################################
336        # Transform questions
337        ################################
338        qaPairs = addQuestions(labeledAnswers, text)
339
340        ################################
341        # Pick the best questions
342        ################################
343        orderedQaPairs = sortAnswers(qaPairs)
344
345        ################################
346        # Generate distractors
347        ################################
348        questions = addDistractors(orderedQaPairs[:count], 4)
349
```
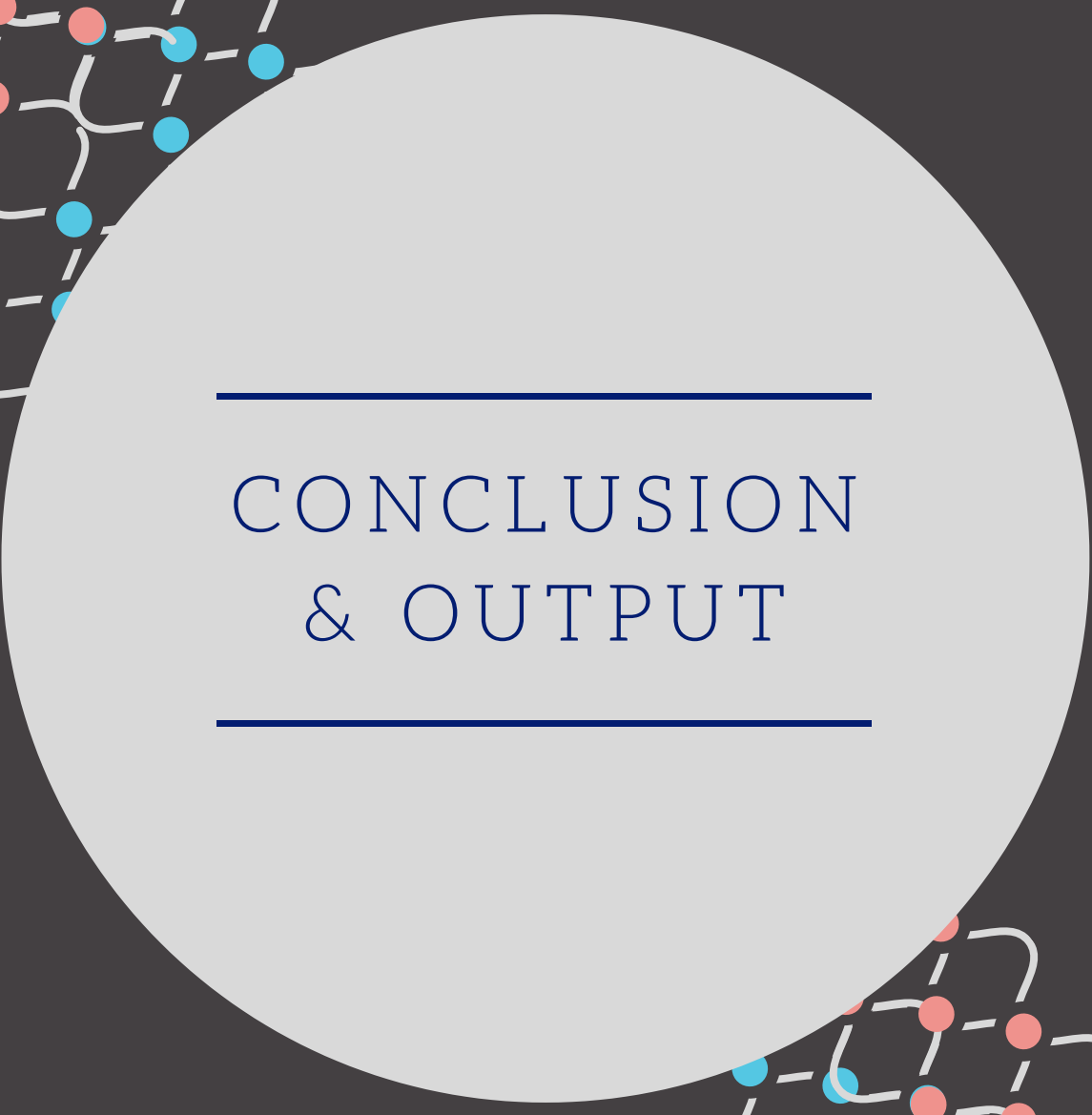
**9**

## STEP 10- VERIFY MODEL OUTCOME AND WRITE MODEL OUTCOME FOR FURTHER ANALYSIS

Yes, This is our MAIN function which actually generates Questions. Here we integrate all the functions to end our Model.

```
350        ################################
351        # Print
352        ################################
353        for i in range(count):
354            options = []
355            options.append(questions[i]['answer'])
356
357            display(Markdown('### Question ' + str(i + 1) + ':'))
358            print(questions[i]['question'])
359
360
361
362            display(Markdown('#### Options:'))
363            for distractor in questions[i]['distractors']:
364                options.append(distractor)
365    #            print(distractor)
366
```

```
369              #################################
370              # Shuffling options
371              #################################
372
373              random.shuffle(options)
374              for num,letter in enumerate(options):
375                  print(num+1," ",letter)
376
377    #          print(ans)
378              display(Markdown('#### Answer:'))
379              for x,correct in enumerate(options):
380                  if correct==questions[i]['answer']:
381                      print(x+1,correct)
382              print()
383    f = open(vAR_Test_Data,mode='r')
384    vAR_Content = f.read()
385    print(vAR_Content)
386    display(Markdown('#### Content'))
387    print('')
388
389    generateQuestions(vAR_Content, 15)
390
391    /*********************************************************************
392    Disclaimer.
393
394    We are providing this code block strictly for learning and researching,this is not a
395    production ready code. We have no liability on this particular code under any circumstances;
396    Users should use this code on their own risk. All software, hardware and other products
397    that are referenced in these materials belong to the respective vendor who developed or who
398    owns this product.
399    /*********************************************************************
400
```

# CONCLUSION & OUTPUT

# Conclusion

We used Gaussian Naïve Bayes Model to Generate Multiple Choice Questions and similar distractors to the answer. The Model Performed well on the test data & predicted the outcome expected.

# Output

**Question 3:**

Furthermore, many experts _____ AI could solve major challenges and crisis situations.

**Options:**

1    believed
2    know
3    believe
4    say
5    think

**Answer:**

3 believe

**Question 4:**

Furthermore, Machine learning algorithms _____ in better serving customers.

**Options:**
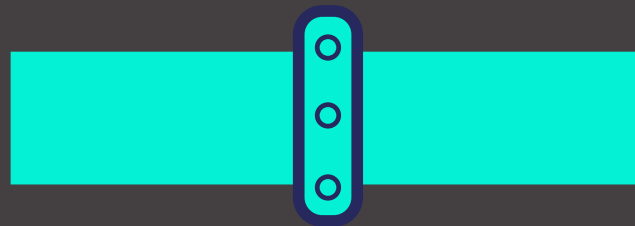
1    helps
2    to
3    helped
4    helping
5    help

**Answer:**

5 help

# APPENDIX

# Kubeflow Pipelines

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Kubeflow Pipelines

> ## Components and Functions

The Kubeflow Pipelines platform consists of:

- A user interface (UI) for managing and tracking experiments, jobs, and runs.
- An engine for scheduling multi-step ML workflows.
- An SDK for defining and manipulating pipelines and components.
- Notebooks for interacting with the system using the SDK.

The Kubeflow Pipelines platform has the following goals:

- End-to-end orchestration: enabling and simplifying the orchestration of machine learning pipelines.
- Easy experimentation: making it easy to try numerous ideas and techniques and manage various trials/experiments.
- Easy re-use: enabling to re-use components and pipelines to quickly cobble together end-to-end solutions, without having to rebuild each time.
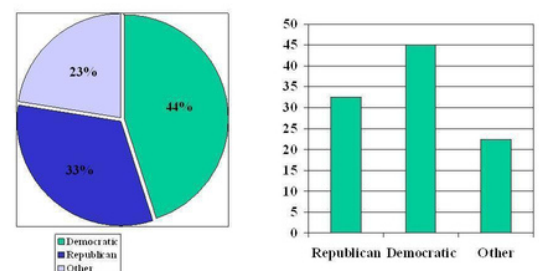
# Types of Data in Machine Learning

DeepSphere.AI
Enterprise AI and IIoT for Analytics

## Types of Data in Machine Learning

# "Nominal Data"

**Nominal values represent discrete units and are used to label variables that have no quantitative value. Just think of them as „labels". Note that nominal data that has no order. Therefore if you would change the order of its values, the meaning would not change. You can see two examples of nominal features below:**

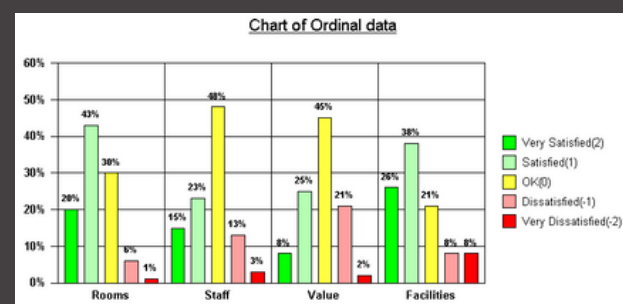Sample Pie Charts and Bar Charts of Nominal Data

Democratic
Republican
Other

23%
44%
33%

50
45
40
35
30
25
20
15
10
5
0

Republican   Democratic   Other

Anthony J Greene                16

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Types of Data in Machine Learning

# " Ordinal Data "

In ordinal encoding, each unique category value is assigned an integer value.For example, "red" is 1, "green" is 2, and "blue" is 3.This is called an ordinal encoding or an integer encoding and is easily reversible. Often, integer values starting at zero are used.For some variables, an ordinal encoding may be enough. The integer values have a natural ordered relationship between each other and machine learning algorithms may be able to understand and harness this relationship.


Chart of Ordinal data

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Types of Data in Machine Learning

" **Continuous Data** "

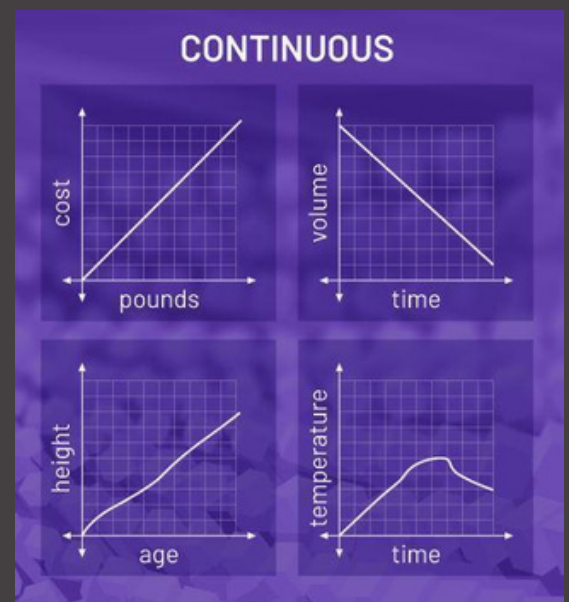In ordinal encoding, each unique category value is assigned an integer value.For example, "red" is 1, "green" is 2, and "blue" is 3.This is called an ordinal encoding or an integer encoding and is easily reversible. Often, integer values starting at zero are used.For some variables, an ordinal encoding may be enough. The integer values have a natural ordered relationship between each other and machine learning algorithms may be able to understand and harness this relationship.

DeepSphere.AI
Enterprise AI and IIoT for Analytics

# Validation in Machine Learning

**Validation in Machine Learning**

> # Validation Actual

In ordinal encoding, each unique category value is assigned an integer value.For example, "red" is 1, "green" is 2, and "blue" is 3.This is called an ordinal encoding or an integer encoding and is easily reversible. Often, integer values starting at zero are used.For some variables, an ordinal encoding may be enough. The integer values have a natural ordered relationship between each other and machine learning algorithms may be able to understand and harness this relationship.
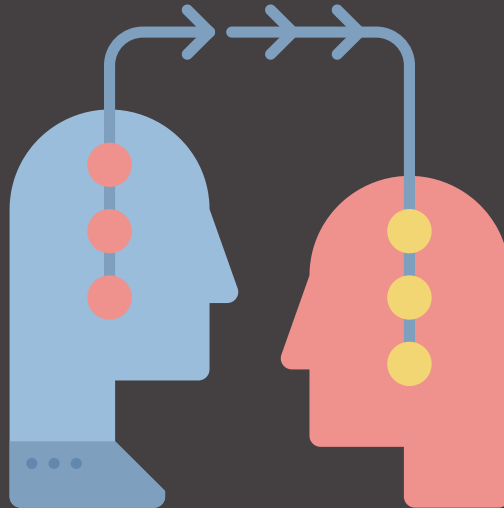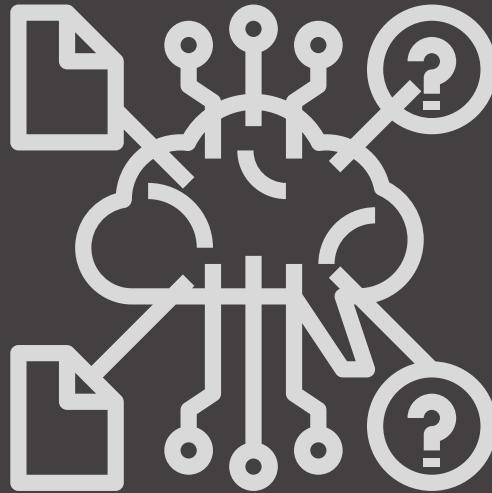
# Training

**Training**

> # Training - Actual

The process of training a DL model involves providing a DL algorithm (that is, the learning algorithm) with training data to learn from. The term DL model refers to the model artifact that is created by the training process. You can use the DL model to get predictions on new data for which you do not know the target.

Training

# " Training - Error "

**Training error is the error that you get when you run the trained model back on the training data. Remember that this data has already been used to train the model and this necessarily doesn't mean that the model once trained will accurately perform when applied back on the training data itself.**

# Prediction

**Prediction**

> # Prediction in Deep Learning

"Prediction" refers to the output of an algorithm after it has been trained on a historical dataset and applied to new data when forecasting the likelihood of a particular outcome, such as whether or not a customer will churn in 30 days. The algorithm will generate probable values for an unknown variable for each record in the new data, allowing the model builder to identify what that value will most likely be.
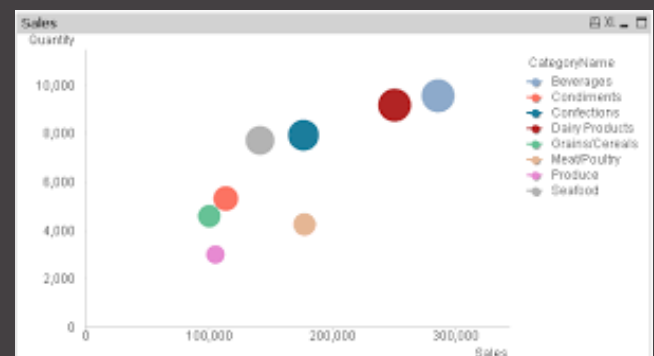
DeepSphere.AI
Enterprise AI and IIoT for Analytics

# Data Visualisations

## Data Visualisations
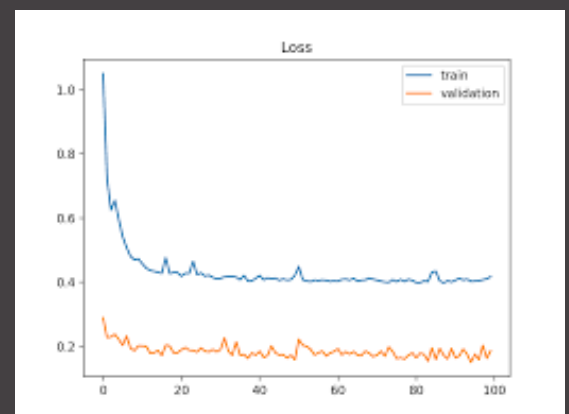
" 

# Bubble Chart

"

A bubble chart is a data visualization that displays multiple circles (bubbles) in a two-dimensional plot. It is a generalization of the scatter plot, replacing the dots with bubbles.

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Data Visualisations

# " Line Chart "

A line chart is, as one can imagine, a line or multiple lines showing how single, or multiple variables develop over time. It is a great tool because we can easily highlight the magnitude of change of one or more variables over a period.
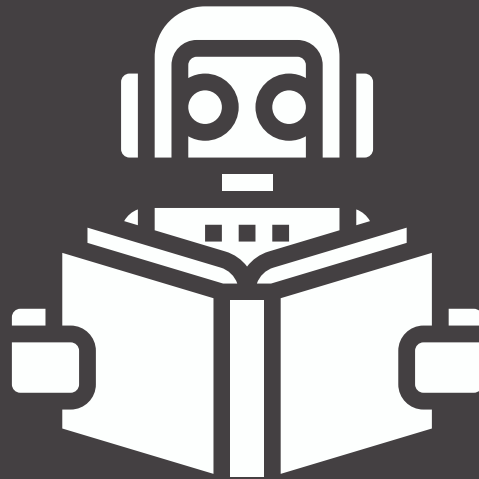
**Data Visualisations**

# "Influencers"

**Influencers are fields that you suspect contain information about someone or something that influences or contributes to anomalies in your data.Influencers can be any field in your data. If you use data feeds, however, the field must exist in your data feed query or aggregation; otherwise it is not included in the job analysis.**

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Data Visualisations

" 

# Standard Deviation

"

**Standard deviation is a number that describes how spread out the values are.A low standard deviation means that most of the numbers are close to the mean (average) value.A high standard deviation means that the values are spread out over a wider range.**

# Some Common ML Terms

**Some Common ML Terms**

"

# Density

"

Use statistical models to find an underlying probability distribution that gives rise to the observed variables.

## Some Common ML Terms

> # Lorenz Curve

**Lorenz curve is also known under the name of "lift curve" when applied to classification/ranking. For a given range of predicted probability values, the lift represents a multiplicative increase in the positive class's rate (due to a given predictive model) over a random guess.**

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Some Common ML Terms

" **Sensitivity** "

**Sensitivity is a measure of the proportion of actual positive cases that got predicted as positive (or true positive). This implies that there will be another proportion of actual positive cases, which would get predicted incorrectly as negative (and, thus, could also be termed as the false negative).**

**DeepSphere.AI**
Enterprise AI and IIoT for Analytics

## Some Common ML Terms

"

# Lift

"

In data mining and association rule learning, lift is a measure of the performance of a targeting model (association rule) at predicting or classifying cases as having an enhanced response (with respect to the population as a whole), measured against a random choice targeting model.