

text_classification_readme

November 13, 2023

1 Text Classification

1.1 Introduction

Emails, chat exchanges, webpages, and social media all contain unstructured content. However, obtaining value from this data is a challenge unless it is structured in a certain way. It used to be a laborious and costly procedure since it needed time and resources to manually analyze the data. Text classifiers utilizing NLP have proved to be an excellent choice for quickly, cost-effectively, and scalable text data structure.

1.2 What is Text Classification?

The technique of categorizing text into structured groupings is known as text classification, alternatively known as text tagging or text categorization. Text classifiers can automatically evaluate text and assign a set of pre-defined tags or categories depending on its content using Natural Language Processing (NLP)

1.3 Some examples of Text Classification

Since it makes it easier to gain insights from data and automate operations inside an organization, text classification is becoming increasingly crucial to industries. Some of the most popular examples and use cases for automated text classification include:

Sentiment Analysis: Analysis of a particular text to determine if it expresses good or negative feelings towards the topic at hand (e.g. for brand monitoring purposes). #### **Topic Detection:** Finding a topic of a line of text (e.g. to know if a product review is about Ease of Use, Customer Support, or Pricing when analyzing customer feedback). #### **Language Detection:** It is the process to determine the language of a text (e.g. know if an incoming support ticket is written in English or Local language for directing the tickets to the appropriate team).

1.4 Understanding the models

Convolutional Neural Network A neural network in which at least one layer is a convolutional layer. A typical convolutional neural network consists of some combination of the following layers:
* convolutional layers * pooling layers * dense layers

Convolutional neural networks have had great success in certain kinds of problems, such as image recognition and text classification.

Recurrent Neural Network A neural network that is intentionally run multiple times, where parts of each run feed into the next run. Specifically, hidden layers from the previous run provide part of the input to the same hidden layer in the next run. Recurrent neural networks are particularly useful for evaluating sequences so that the hidden layers can learn from previous runs of the neural network on earlier parts of the sequence.

Long Short-Term Memory (LSTM) A type of cell in a recurrent neural network is used to process sequences of data in applications such as handwriting recognition, machine translation, and image captioning. LSTMs address the vanishing gradient problem that occurs when training RNNs due to long data sequences by maintaining history in an internal memory state based on new input and context from previous cells in the RNN.

1.5 Why this approach?

When we use CNN on text data through a diagram, the result of each convolution will fire when a special pattern is detected. By varying the size of the kernels and concatenating their outputs, you're allowing yourself to detect patterns of multiples sizes (2, 3, or 5 adjacent words). Patterns could be expressions (word ngrams?) like "I hate", "very good" and therefore CNNs can identify them in the sentence regardless of their position. Recurrent neural networks can obtain context information but the order of words will lead to bias; the text analysis method based on Convolutional neural network (CNN) can obtain important features of text through pooling but it is difficult to obtain contextual information which can be leverage using LSTM. So using the combination of CNN with LSTM could give us some interesting results

1.6 A sample problem we are trying to solve

Article review sentiment classification problem. Each review is a variable sequence of words and the sentiment of each review must be classified. The Article Review Dataset contains 999 highly-polar reviews (good or bad) stored as text files. The 699 text files for training and the 299 for testing. The problem is to determine whether a given article review has a positive or negative sentiment.

1.7 Data Exploration and Pre-processing

First, I use binary encoding for the sentiments, i.e $y = 1$ for positive sentiments and $y = -1$ for negative sentiments. Since the provided data are pretty clean, we can remove the punctuation and numbers from the data for further analysis.

```
In [4]: #pos = os.listdir("../data/pos")
#neg = os.listdir("../data/neg")

from google.colab import drive
drive.mount('/content/drive')
drive_folder = '/content/drive/My Drive/552project/data'

pos = os.listdir(drive_folder + '/pos/')
neg = os.listdir(drive_folder + '/neg/')

data = []
punc_dig = string.punctuation + string.digits

for path in pos:
    f = open("/content/drive/My Drive/552project/data/pos/" + path)
    review = f.readlines()
    f.close()
    for i in range(len(review)):
        review[i] = review[i].translate(str.maketrans('', '', punc_dig))
    data.append([int(path[2:5]), " ".join(review), 1])

for path in neg:
    f = open("/content/drive/My Drive/552project/data/neg/" + path)
    review = f.readlines()
    f.close()
    for i in range(len(review)):
        review[i] = review[i].translate(str.maketrans('', '', punc_dig))
    data.append([int(path[2:5]), " ".join(review), -1])
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

The name of each text file starts with cv number. Here, I split the data use text files 0-699 in each class for training and 700-999 for testing.

```
In [5]: data = pd.DataFrame(data, columns = ["num", "text", "class"])
```

```
In [6]: train = data[data["num"] < 700]
test = data[data["num"] >= 700]
```

Count the number of unique words in the whole dataset (train + test) and print it out.

```
In [7]: # resource: https://stackoverflow.com/questions/18936957/count-distinct-words-from-a-pandas-data-frame

results = set()
data['text'].str.lower().str.split().apply(results.update)
print("number of unique words: ", len(results))

number of unique words: 46830
```

```
In [ ]: #result = data.text.apply(lambda x: pd.value_counts(x.split(" ")).sum(axis = 0))
```

Calculate the average review length and the standard deviation of review lengths. Report the results.

```
In [ ]: #data0 = data.replace("\n", " ")
```

```
In [8]: review_len = []  
  
for i in data["text"]:  
    review_len.append(len(i.split()))  
  
print("average: ", mean(review_len))  
print("standard deviation: ", stdev(review_len))
```

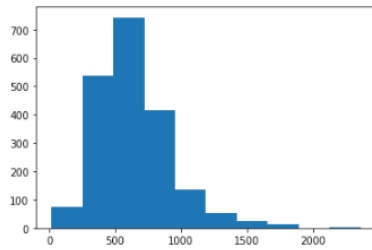
```
average: 644.3555  
standard deviation: 285.0511431249635
```

Visualization

Plot the histogram of review lengths.

```
In [ ]: plt.hist(review_len)
```

```
Out[10]: (array([ 75., 537., 743., 417., 136., 53., 23., 13., 0., 3.]),  
array([ 16., 250.7, 485.4, 720.1, 954.8, 1189.5, 1424.2, 1658.9,  
1893.6, 2128.3, 2363. ]),  
<BarContainer object of 10 artists>)
```



1.8 Word Embeddings

NLP/Deep Learning

To represent each text (= data point), there are many ways. In NLP/Deep Learning terminology, this task is called tokenization. It is common to represent text using popularity/ rank of words in text. The most common word in the text will be represented as 1, the second most common word will be represented as 2, etc. Tokenize each text document using this method.

```
In [ ]: token = Tokenizer()
token.fit_on_texts(data["text"])
pop = token.word_index
pop

{'him': 54,
 'into': 55,
 'even': 56,
 'only': 57,
 'than': 58,
 'no': 59,
 'we': 60,
 'good': 61,
 'most': 62,
 'time': 63,
 'can': 64,
 'will': 65,
 'story': 66,
 'films': 67,
 'been': 68,
 'would': 69,
 'much': 70,
 'also': 71,
 'characters': 72,
 'text': 73,
 'the': 74,
 'a': 75,
 'and': 76,
 'of': 77,
 'is': 78,
 'on': 79,
 'to': 80,
 'in': 81,
 'at': 82,
 'from': 83,
 'with': 84,
 'by': 85,
 'for': 86,
 'as': 87,
 'but': 88,
 'or': 89,
 'so': 90,
 'that': 91,
 'this': 92,
 'which': 93,
 'what': 94,
 'who': 95,
 'how': 96,
 'when': 97,
 'where': 98,
 'why': 99,
 'if': 100,
 'then': 101,
 'and': 102,
 'or': 103,
 'but': 104,
 'so': 105,
 'that': 106,
 'this': 107,
 'which': 108,
 'what': 109,
 'who': 110,
 'how': 111,
 'when': 112,
 'where': 113,
 'why': 114,
 'if': 115,
 'then': 116,
 'and': 117,
 'or': 118,
 'but': 119,
 'so': 120,
 'that': 121,
 'this': 122,
 'which': 123,
 'what': 124,
 'who': 125,
 'how': 126,
 'when': 127,
 'where': 128,
 'why': 129,
 'if': 130,
 'then': 131,
 'and': 132,
 'or': 133,
 'but': 134,
 'so': 135,
 'that': 136,
 'this': 137,
 'which': 138,
 'what': 139,
 'who': 140,
 'how': 141,
 'when': 142,
 'where': 143,
 'why': 144,
 'if': 145,
 'then': 146,
 'and': 147,
 'or': 148,
 'but': 149,
 'so': 150,
 'that': 151,
 'this': 152,
 'which': 153,
 'what': 154,
 'who': 155,
 'how': 156,
 'when': 157,
 'where': 158,
 'why': 159,
 'if': 160,
 'then': 161,
 'and': 162,
 'or': 163,
 'but': 164,
 'so': 165,
 'that': 166,
 'this': 167,
 'which': 168,
 'what': 169,
 'who': 170,
 'how': 171,
 'when': 172,
 'where': 173,
 'why': 174,
 'if': 175,
 'then': 176,
 'and': 177,
 'or': 178,
 'but': 179,
 'so': 180,
 'that': 181,
 'this': 182,
 'which': 183,
 'what': 184,
 'who': 185,
 'how': 186,
 'when': 187,
 'where': 188,
 'why': 189,
 'if': 190,
 'then': 191,
 'and': 192,
 'or': 193,
 'but': 194,
 'so': 195,
 'that': 196,
 'this': 197,
 'which': 198,
 'what': 199,
 'who': 200,
 'how': 201,
 'when': 202,
 'where': 203,
 'why': 204,
 'if': 205,
 'then': 206,
 'and': 207,
 'or': 208,
 'but': 209,
 'so': 210,
 'that': 211,
 'this': 212,
 'which': 213,
 'what': 214,
 'who': 215,
 'how': 216,
 'when': 217,
 'where': 218,
 'why': 219,
 'if': 220,
 'then': 221,
 'and': 222,
 'or': 223,
 'but': 224,
 'so': 225,
 'that': 226,
 'this': 227,
 'which': 228,
 'what': 229,
 'who': 230,
 'how': 231,
 'when': 232,
 'where': 233,
 'why': 234,
 'if': 235,
 'then': 236,
 'and': 237,
 'or': 238,
 'but': 239,
 'so': 240,
 'that': 241,
 'this': 242,
 'which': 243,
 'what': 244,
 'who': 245,
 'how': 246,
 'when': 247,
 'where': 248,
 'why': 249,
 'if': 250,
 'then': 251,
 'and': 252,
 'or': 253,
 'but': 254,
 'so': 255,
 'that': 256,
 'this': 257,
 'which': 258,
 'what': 259,
 'who': 260,
 'how': 261,
 'when': 262,
 'where': 263,
 'why': 264,
 'if': 265,
 'then': 266,
 'and': 267,
 'or': 268,
 'but': 269,
 'so': 270,
 'that': 271,
 'this': 272,
 'which': 273,
 'what': 274,
 'who': 275,
 'how': 276,
 'when': 277,
 'where': 278,
 'why': 279,
 'if': 280,
 'then': 281,
 'and': 282,
 'or': 283,
 'but': 284,
 'so': 285,
 'that': 286,
 'this': 287,
 'which': 288,
 'what': 289,
 'who': 290,
 'how': 291,
 'when': 292,
 'where': 293,
 'why': 294,
 'if': 295,
 'then': 296,
 'and': 297,
 'or': 298,
 'but': 299,
 'so': 300,
 'that': 301,
 'this': 302,
 'which': 303,
 'what': 304,
 'who': 305,
 'how': 306,
 'when': 307,
 'where': 308,
 'why': 309,
 'if': 310,
 'then': 311,
 'and': 312,
 'or': 313,
 'but': 314,
 'so': 315,
 'that': 316,
 'this': 317,
 'which': 318,
 'what': 319,
 'who': 320,
 'how': 321,
 'when': 322,
 'where': 323,
 'why': 324,
 'if': 325,
 'then': 326,
 'and': 327,
 'or': 328,
 'but': 329,
 'so': 330,
 'that': 331,
 'this': 332,
 'which': 333,
 'what': 334,
 'who': 335,
 'how': 336,
 'when': 337,
 'where': 338,
 'why': 339,
 'if': 340,
 'then': 341,
 'and': 342,
 'or': 343,
 'but': 344,
 'so': 345,
 'that': 346,
 'this': 347,
 'which': 348,
 'what': 349,
 'who': 350,
 'how': 351,
 'when': 352,
 'where': 353,
 'why': 354,
 'if': 355,
 'then': 356,
 'and': 357,
 'or': 358,
 'but': 359,
 'so': 360,
 'that': 361,
 'this': 362,
 'which': 363,
 'what': 364,
 'who': 365,
 'how': 366,
 'when': 367,
 'where': 368,
 'why': 369,
 'if': 370,
 'then': 371,
 'and': 372,
 'or': 373,
 'but': 374,
 'so': 375,
 'that': 376,
 'this': 377,
 'which': 378,
 'what': 379,
 'who': 380,
 'how': 381,
 'when': 382,
 'where': 383,
 'why': 384,
 'if': 385,
 'then': 386,
 'and': 387,
 'or': 388,
 'but': 389,
 'so': 390,
 'that': 391,
 'this': 392,
 'which': 393,
 'what': 394,
 'who': 395,
 'how': 396,
 'when': 397,
 'where': 398,
 'why': 399,
 'if': 400,
 'then': 401,
 'and': 402,
 'or': 403,
 'but': 404,
 'so': 405,
 'that': 406,
 'this': 407,
 'which': 408,
 'what': 409,
 'who': 410,
 'how': 411,
 'when': 412,
 'where': 413,
 'why': 414,
 'if': 415,
 'then': 416,
 'and': 417,
 'or': 418,
 'but': 419,
 'so': 420,
 'that': 421,
 'this': 422,
 'which': 423,
 'what': 424,
 'who': 425,
 'how': 426,
 'when': 427,
 'where': 428,
 'why': 429,
 'if': 430,
 'then': 431,
 'and': 432,
 'or': 433,
 'but': 434,
 'so': 435,
 'that': 436,
 'this': 437,
 'which': 438,
 'what': 439,
 'who': 440,
 'how': 441,
 'when': 442,
 'where': 443,
 'why': 444,
 'if': 445,
 'then': 446,
 'and': 447,
 'or': 448,
 'but': 449,
 'so': 450,
 'that': 451,
 'this': 452,
 'which': 453,
 'what': 454,
 'who': 455,
 'how': 456,
 'when': 457,
 'where': 458,
 'why': 459,
 'if': 460,
 'then': 461,
 'and': 462,
 'or': 463,
 'but': 464,
 'so': 465,
 'that': 466,
 'this': 467,
 'which': 468,
 'what': 469,
 'who': 470,
 'how': 471,
 'when': 472,
 'where': 473,
 'why': 474,
 'if': 475,
 'then': 476,
 'and': 477,
 'or': 478,
 'but': 479,
 'so': 480,
 'that': 481,
 'this': 482,
 'which': 483,
 'what': 484,
 'who': 485,
 'how': 486,
 'when': 487,
 'where': 488,
 'why': 489,
 'if': 490,
 'then': 491,
 'and': 492,
 'or': 493,
 'but': 494,
 'so': 495,
 'that': 496,
 'this': 497,
 'which': 498,
 'what': 499,
 'who': 500,
 'how': 501,
 'when': 502,
 'where': 503,
 'why': 504,
 'if': 505,
 'then': 506,
 'and': 507,
 'or': 508,
 'but': 509,
 'so': 510,
 'that': 511,
 'this': 512,
 'which': 513,
 'what': 514,
 'who': 515,
 'how': 516,
 'when': 517,
 'where': 518,
 'why': 519,
 'if': 520,
 'then': 521,
 'and': 522,
 'or': 523,
 'but': 524,
 'so': 525,
 'that': 526,
 'this': 527,
 'which': 528,
 'what': 529,
 'who': 530,
 'how': 531,
 'when': 532,
 'where': 533,
 'why': 534,
 'if': 535,
 'then': 536,
 'and': 537,
 'or': 538,
 'but': 539,
 'so': 540,
 'that': 541,
 'this': 542,
 'which': 543,
 'what': 544,
 'who': 545,
 'how': 546,
 'when': 547,
 'where': 548,
 'why': 549,
 'if': 550,
 'then': 551,
 'and': 552,
 'or': 553,
 'but': 554,
 'so': 555,
 'that': 556,
 'this': 557,
 'which': 558,
 'what': 559,
 'who': 560,
 'how': 561,
 'when': 562,
 'where': 563,
 'why': 564,
 'if': 565,
 'then': 566,
 'and': 567,
 'or': 568,
 'but': 569,
 'so': 570,
 'that': 571,
 'this': 572,
 'which': 573,
 'what': 574,
 'who': 575,
 'how': 576,
 'when': 577,
 'where': 578,
 'why': 579,
 'if': 580,
 'then': 581,
 'and': 582,
 'or': 583,
 'but': 584,
 'so': 585,
 'that': 586,
 'this': 587,
 'which': 588,
 'what': 589,
 'who': 590,
 'how': 591,
 'when': 592,
 'where': 593,
 'why': 594,
 'if': 595,
 'then': 596,
 'and': 597,
 'or': 598,
 'but': 599,
 'so': 600,
 'that': 601,
 'this': 602,
 'which': 603,
 'what': 604,
 'who': 605,
 'how': 606,
 'when': 607,
 'where': 608,
 'why': 609,
 'if': 610,
 'then': 611,
 'and': 612,
 'or': 613,
 'but': 614,
 'so': 615,
 'that': 616,
 'this': 617,
 'which': 618,
 'what': 619,
 'who': 620,
 'how': 621,
 'when': 622,
 'where': 623,
 'why': 624,
 'if': 625,
 'then': 626,
 'and': 627,
 'or': 628,
 'but': 629,
 'so': 630,
 'that': 631,
 'this': 632,
 'which': 633,
 'what': 634,
 'who': 635,
 'how': 636,
 'when': 637,
 'where': 638,
 'why': 639,
 'if': 640,
 'then': 641,
 'and': 642,
 'or': 643,
 'but': 644,
 'so': 645,
 'that': 646,
 'this': 647,
 'which': 648,
 'what': 649,
 'who': 650,
 'how': 651,
 'when': 652,
 'where': 653,
 'why': 654,
 'if': 655,
 'then': 656,
 'and': 657,
 'or': 658,
 'but': 659,
 'so': 660,
 'that': 661,
 'this': 662,
 'which': 663,
 'what': 664,
 'who': 665,
 'how': 666,
 'when': 667,
 'where': 668,
 'why': 669,
 'if': 670,
 'then': 671,
 'and': 672,
 'or': 673,
 'but': 674,
 'so': 675,
 'that': 676,
 'this': 677,
 'which': 678,
 'what': 679,
 'who': 680,
 'how': 681,
 'when': 682,
 'where': 683,
 'why': 684,
 'if': 685,
 'then': 686,
 'and': 687,
 'or': 688,
 'but': 689,
 'so': 690,
 'that': 691,
 'this': 692,
 'which': 693,
 'what': 694,
 'who': 695,
 'how': 696,
 'when': 697,
 'where': 698,
 'why': 699,
 'if': 700,
 'then': 701,
 'and': 702,
 'or': 703,
 'but': 704,
 'so': 705,
 'that': 706,
 'this': 707,
 'which': 708,
 'what': 709,
 'who': 710,
 'how': 711,
 'when': 712,
 'where': 713,
 'why': 714,
 'if': 715,
 'then': 716,
 'and': 717,
 'or': 718,
 'but': 719,
 'so': 720,
 'that': 721,
 'this': 722,
 'which': 723,
 'what': 724,
 'who': 725,
 'how': 726,
 'when': 727,
 'where': 728,
 'why': 729,
 'if': 730,
 'then': 731,
 'and': 732,
 'or': 733,
 'but': 734,
 'so': 735,
 'that': 736,
 'this': 737,
 'which': 738,
 'what': 739,
 'who': 740,
 'how': 741,
 'when': 742,
 'where': 743,
 'why': 744,
 'if': 745,
 'then': 746,
 'and': 747,
 'or': 748,
 'but': 749,
 'so': 750,
 'that': 751,
 'this': 752,
 'which': 753,
 'what': 754,
 'who': 755,
 'how': 756,
 'when': 757,
 'where': 758,
 'why': 759,
 'if': 760,
 'then': 761,
 'and': 762,
 'or': 763,
 'but': 764,
 'so': 765,
 'that': 766,
 'this': 767,
 'which': 768,
 'what': 769,
 'who': 770,
 'how': 771,
 'when': 772,
 'where': 773,
 'why': 774,
 'if': 775,
 'then': 776,
 'and': 777,
 'or': 778,
 'but': 779,
 'so': 780,
 'that': 781,
 'this': 782,
 'which': 783,
 'what': 784,
 'who': 785,
 'how': 786,
 'when': 787,
 'where': 788,
 'why': 789,
 'if': 790,
 'then': 791,
 'and': 792,
 'or': 793,
 'but': 794,
 'so': 795,
 'that': 796,
 'this': 797,
 'which': 798,
 'what': 799,
 'who': 800,
 'how': 801,
 'when': 802,
 'where': 803,
 'why': 804,
 'if': 805,
 'then': 806,
 'and': 807,
 'or': 808,
 'but': 809,
 'so': 810,
 'that': 811,
 'this': 812,
 'which': 813,
 'what': 814,
 'who': 815,
 'how': 816,
 'when': 817,
 'where': 818,
 'why': 819,
 'if': 820,
 'then': 821,
 'and': 822,
 'or': 823,
 'but': 824,
 'so': 825,
 'that': 826,
 'this': 827,
 'which': 828,
 'what': 829,
 'who': 830,
 'how': 831,
 'when': 832,
 'where': 833,
 'why': 834,
 'if': 835,
 'then': 836,
 'and': 837,
 'or': 838,
 'but': 839,
 'so': 840,
 'that': 841,
 'this': 842,
 'which': 843,
 'what': 844,
 'who': 845,
 'how': 846,
 'when': 847,
 'where': 848,
 'why': 849,
 'if': 850,
 'then': 851,
 'and': 852,
 'or': 853,
 'but': 854,
 'so': 855,
 'that': 856,
 'this': 857,
 'which': 858,
 'what': 859,
 'who': 860,
 'how': 861,
 'when': 862,
 'where': 863,
 'why': 864,
 'if': 865,
 'then': 866,
 'and': 867,
 'or': 868,
 'but': 869,
 'so': 870,
 'that': 871,
 'this': 872,
 'which': 873,
 'what': 874,
 'who': 875,
 'how': 876,
 'when': 877,
 'where': 878,
 'why': 879,
 'if': 880,
 'then': 881,
 'and': 882,
 'or': 883,
 'but': 884,
 'so': 885,
 'that': 886,
 'this': 887,
 'which': 888,
 'what': 889,
 'who': 890,
 'how': 891,
 'when': 892,
 'where': 893,
 'why': 894,
 'if': 895,
 'then': 896,
 'and': 897,
 'or': 898,
 'but': 899,
 'so': 900,
 'that': 901,
 'this': 902,
 'which': 903,
 'what': 904,
 'who': 905,
 'how': 906,
 'when': 907,
 'where': 908,
 'why': 909,
 'if': 910,
 'then': 911,
 'and': 912,
 'or': 913,
 'but': 914,
 'so': 915,
 'that': 916,
 'this': 917,
 'which': 918,
 'what': 919,
 'who': 920,
 'how': 921,
 'when': 922,
 'where': 923,
 'why': 924,
 'if': 925,
 'then': 926,
 'and': 927,
 'or': 928,
 'but': 929,
 'so': 930,
 'that': 931,
 'this': 932,
 'which': 933,
 'what': 934,
 'who': 935,
 'how': 936,
 'when': 937,
 'where': 938,
 'why': 939,
 'if': 940,
 'then': 941,
 'and': 942,
 'or': 943,
 'but': 944,
 'so': 945,
 'that': 946,
 'this': 947,
 'which': 948,
 'what': 949,
 'who': 950,
 'how': 951,
 'when': 952,
 'where': 953,
 'why': 954,
 'if': 955,
 'then': 956,
 'and': 957,
 'or': 958,
 'but': 959,
 'so': 960,
 'that': 961,
 'this': 962,
 'which': 963,
 'what': 964,
 'who': 965,
 'how': 966,
 'when': 967,
 'where': 968,
 'why': 969,
 'if': 970,
 'then': 971,
 'and': 972,
 'or': 973,
 'but': 974,
 'so': 975,
 'that': 976,
 'this': 977,
 'which': 978,
 'what': 979,
 'who': 980,
 'how': 981,
 'when': 982,
 'where': 983,
 'why': 984,
 'if': 985,
 'then': 986,
 'and': 987,
 'or': 988,
 'but': 989,
 'so': 990,
 'that': 991,
 'this': 992,
 'which': 993,
 'what': 994,
 'who': 995,
 'how': 996,
 'when': 997,
 'where': 998,
 'why': 999,
 'if': 1000,
 'then': 1001,
 'and': 1002,
 'or': 1003,
 'but': 1004,
 'so': 1005,
 'that': 1006,
 'this': 1007,
 'which': 1008,
 'what': 1009,
 'who': 1010,
 'how': 1011,
 'when': 1012,
 'where': 1013,
 'why': 1014,
 'if': 1015,
 'then': 1016,
 'and': 1017,
 'or': 1018,
 'but': 1019,
 'so': 1020,
 'that': 1021,
 'this': 1022,
 'which': 1023,
 'what': 1024,
 'who': 1025,
 'how': 1026,
 'when': 1027,
 'where': 1028,
 'why': 1029,
 'if': 1030,
 'then': 1031,
 'and': 1032,
 'or': 1033,
 'but': 1034,
 'so': 1035,
 'that': 1036,
 'this': 1037,
 'which': 1038,
 'what': 1039,
 'who': 1040,
 'how': 1041,
 'when': 1042,
 'where': 1043,
 'why': 1044,
 'if': 1045,
 'then': 1046,
 'and': 1047,
 'or': 1048,
 'but': 1049,
 'so': 1050,
 'that': 1051,
 'this': 1052,
 'which': 1053,
 'what': 1054,
 'who': 1055,
 'how': 1056,
 'when': 1057,
 'where': 1058,
 'why': 1059,
 'if': 1060,
 'then': 1061,
 'and': 1062,
 'or': 1063,
 'but': 1064,
 'so': 1065,
 'that': 1066,
 'this': 1067,
 'which': 1068,
 'what': 1069,
 'who': 1070,
 'how': 1071,
 'when': 1072,
 'where': 1073,
 'why': 1074,
 'if': 1075,
 'then': 1076,
 'and': 1077,
 'or': 1078,
 'but': 1079,
 'so': 1080,
 'that': 1081,
 'this': 1082,
 'which': 1083,
 'what': 1084,
 'who': 1085,
 'how': 1086,
 'when': 1087,
 'where': 1088,
 'why': 1089,
 'if': 1090,
 'then': 1091,
 'and': 1092,
 'or': 1093,
 'but': 1094,
 'so': 1095,
 'that': 1096,
 'this': 1097,
 'which': 1098,
 'what': 1099,
 'who': 1100,
 'how': 1101,
 'when': 1102,
 'where': 1103,
 'why': 1104,
 'if': 1105,
 'then': 1106,
 'and': 1107,
 'or': 1108,
 'but': 1109,
 'so': 1110,
 'that': 1111,
 'this': 1112,
 'which': 1113,
 'what': 1114,
 'who': 1115,
 'how': 1116,
 'when': 1117,
 'where': 1118,
 'why': 1119,
 'if': 1120,
 'then': 1121,
 'and': 1122,
 'or': 1123,
 'but': 1124,
 'so': 1125,
 'that': 1126,
 'this': 1127,
 'which': 1128,
 'what': 1129,
 'who': 1130,
 'how': 1131,
 'when': 1132,
 'where': 1133,
 'why': 1134,
 'if': 1135,
 'then': 1136,
 'and': 1137,
 'or': 1138,
 'but': 1139,
 'so': 1140,
 'that': 1141,
 'this': 1142,
 'which': 1143,
 'what': 1144,
 'who': 1145,
 'how': 1146,
 'when': 1147,
 'where': 1148,
 'why': 1149,
 'if': 1150,
 'then': 1151,
 'and': 1152,
 'or': 1153,
 'but': 1154,
 'so': 1155,
 'that': 1156,
 'this': 1157,
 'which': 1158,
 'what': 1159,
 'who': 1160,
 'how': 1161,
 'when': 1162,
 'where': 1163,
 'why': 1164,
 'if': 1165,
 'then': 1166,
 'and': 1167,
 'or': 1168,
 'but': 1169,
 'so': 1170,
 'that': 1171,
 'this': 1172,
 'which': 1173,
 'what': 1174,
 'who': 1175,
 'how': 1176,
 'when': 1177,
 'where': 1178,
 'why': 1179,
 'if': 1180,
 'then': 1181,
 'and': 1182,
 'or': 1183,
 'but': 1184,
 'so': 1185,
 'that': 1186,
 'this': 1187,
 'which': 1188,
 'what': 1189,
 'who': 1190,
 'how': 1191,
 'when': 1192,
 'where': 1193,
 'why': 1194,
 'if': 1195,
 'then': 1196,
 'and': 1197,
 'or': 1198,
 'but': 1199,
 'so': 1200,
 'that': 1201,
 'this': 1202,
 'which': 1203,
 'what': 1204,
 'who': 1205,
 'how': 1206,
 'when': 1207,
 'where': 1208,
 'why': 1209,
 'if': 1210,
 'then': 1211,
 'and': 1212,
 'or': 1213,
 'but': 1214,
 'so': 1215,
 'that': 1216,
 'this': 1217,
 'which': 1218,
 'what': 1219,
 'who': 1220,
 'how': 1221,
 'when': 1222,
 'where': 1223,
 'why': 1224,
 'if': 1225,
 'then': 1226,
 'and': 1227,
 'or': 1228,
 'but': 1229,
 'so': 1230,
 'that': 1231,
 'this': 1232,
 'which': 1233,
 'what': 1234,
 'who': 1235,
 'how': 1236,
 'when': 1237,
 'where': 1238,
 'why': 1239,
 'if': 1240,
 'then': 1241,
 'and': 1242,
 'or': 1243,
 'but': 1244,
 'so': 1245,
 'that': 1246,
 'this': 1247,
 'which': 1248,
 'what': 1249,
 'who': 1250,
 'how': 1251,
 'when': 1252,
 'where': 1253,
 'why': 1254,
 'if': 1255,
 'then': 1256,
 'and': 1257,
 'or': 1258,
 'but': 1259,
 'so': 1260,
 'that': 1261,
 'this': 1262,
 'which': 1263,
 'what': 1264,
 'who': 1265,
 'how': 1266,
 'when': 1267,
 'where': 1268,
 'why': 1269,
 'if': 1270,
 'then': 1271,
 'and': 1272,
 'or': 1273,
 'but': 1274,
 'so': 1275,
 'that': 1276,
 'this': 1277,
 'which': 1278,
 'what': 1279,
 'who': 1280,
 'how': 1281,
 'when': 1282,
 'where': 1283,
 'why': 1284,
 'if': 1285,
 'then': 1286,
 'and': 1287,
 'or': 1288,
 'but': 1289,
 'so': 1290,
 'that': 1291,
 'this': 1292,
 'which': 1293,
 'what': 1294,
 'who': 1295,
 'how': 1296,
 'when': 1297,
 'where': 1298,
 'why': 1299,
 'if': 1300,
 'then': 1301,
 'and': 1302,
 'or': 1303,
 'but': 1304,
 'so': 1305,
 'that': 1306,
 'this': 1307,
 'which': 1308,
 'what': 1309,
 'who': 1310,
 'how': 1311,
 'when': 1312,
 'where': 1313,
 'why': 1314,
 'if': 1315,
 'then': 1316,
 'and': 1317,
 'or': 1318,
 'but': 1319,
 'so': 1320,
 'that': 1321,
 'this': 1322,
 'which': 1323,
 'what': 1324,
 'who': 1325,
 'how': 1326,
 'when': 1327,
 'where': 1328,
 'why': 1329,
 'if': 1330,
 'then': 1331,
 'and': 1332,
 'or': 1333,
 'but': 1334,
 'so': 1335,
 'that': 1336,
 'this': 1337,
 'which': 1338,
 'what': 1339,
 'who': 1340,
 'how': 1341,
 'when': 1342,
 'where': 1343,
 'why': 1344,
 'if': 1345,
 'then': 1346,
 'and': 1347,
 'or': 1348,
 'but': 1349,
 'so': 1350,
 'that': 1351,
 'this': 1352,
 'which': 1353,
 'what': 1354,
 'who': 1355,
 'how': 1356,
 'when': 1357,
 'where': 1358,
 'why': 1359,
 'if': 1360,
 'then': 1361,
 'and': 1362,
 'or': 1363,
 'but': 1364,
 'so': 1365,
 'that': 1366
```

Word Embeddings

i. Use tokenized text as inputs to a deep neural network.

However, a recent breakthrough in NLP suggests that more sophisticated representations of text yield better results. These sophisticated representations are called word embeddings. Word embedding is a term used for representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning.⁴ Most deep learning modules (including Keras) provide a convenient way to convert positive integer representations of words into a word embedding by an 'Embedding layer.' The layer accepts arguments that define the mapping of words into embeddings, including the maximum number of expected words also called the vocabulary size (e.g. the largest integer value). The layer also allows you to specify the dimension for each word vector, called the 'output dimension.' We would like to use a word embedding layer for this project. Assume that we are interested in the top 5,000 words. This means that in each integer sequence that represents each document, we set to zero those integers that represent words that are not among the top 5,000 words in the document.⁵ If you feel more adventurous, use all the words that appear in this corpus. Choose the length of the embedding vector for each word to be 32. Hence, each document is represented as a $32 \times L$ matrix.

ii. Flatten the matrix of each document to a vector.

```
In [ ]: input_text = np.array(list(data["pad_text"]))
        input_text[input_text >= 5000] = 0

In [ ]: #resource: https://www.cnblogs.com/Renyi-Fan/p/13809918.html
        model = Sequential()
        model.add(Embedding(5000, 32, input_length = L))
        model.add(Flatten())
        model.compile('rmsprop', 'mse')
        print(model.summary)

        embed_text = model.predict(input_text)
        print(model.layers[0].get_weights())
        print(embed_text)

<bound method Model.summary of <tensorflow.python.keras.engine.sequential.Sequential object at 0x160e69df0>>
[array([[[-0.03551153, -0.0312135, 0.02762629, ..., -0.00306189,
         0.03902471, -0.01229275],
        [ 0.04268488, -0.04910192, -0.04501715, ..., 0.00510366,
         0.04548797, -0.00195887],
        [ 0.03849134, 0.03062275, 0.03373922, ..., 0.01225618,
         -0.02856541, 0.0476712 ],
        ...,
        [ 0.02141747, 0.04477562, -0.00798661, ..., -0.0272519,
         0.02140928, 0.04162479],
        [ 0.03989843, 0.04823219, -0.00834541, ..., -0.02606398,
         -0.00012935, -0.02234277],
        [-0.01735542, 0.02892314, 0.02645124, ..., -0.00246954,
         -0.02422092, 0.02956215]], dtype=float32)]
[[[ 0.03333518 -0.02617295 -0.02261096 ... -0.04879908 -0.02016357
    -0.00410762]
  [-0.03551153 -0.0312135 0.02762629 ... -0.00327522 -0.00395964
    0.02752754]
  [-0.04900008 -0.04394411 -0.0189552 ... -0.00306189 0.03902471
    -0.01229275]
  ...
  [-0.03551153 -0.0312135 0.02762629 ... 0.0123669 0.03612501
    -0.03114016]
  [-0.03551153 -0.0312135 0.02762629 ... -0.01558664 0.01954446
    -0.01011112]
  [-0.03551153 -0.0312135 0.02762629 ... -0.04458426 0.02998788
    -0.02043418]]]
```

1.9 Multi-Layer Perceptron

Multi-Layer Perceptron

```
In [ ]: Data = data[["num", "class"]]
Data["emb_text"] = list(embed_text)
Data["class"] = Data["class"].replace(-1, 0)

<ipython-input-18-5e7af8508837>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
Data["emb_text"] = list(embed_text)
<ipython-input-18-5e7af8508837>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
Data["class"] = Data["class"].replace(-1, 0)

In [ ]: Train = Data[Data["num"] < 700]
Test = Data[Data["num"] >= 700]
X_tr = np.array(list(Train["emb_text"]))
X_t = np.array(list(Test["emb_text"]))
Y_tr = Train["class"]
Y_t = Test["class"]
```

i. Train a MLP with three (dense) hidden layers each of which has 50 ReLUs and one output layer with a single sigmoid neuron. Use a dropout rate of 20% for the first layer and 50% for the other layers. Use ADAM optimizer and binary cross entropy loss (which is equivalent to having a softmax in the output). To avoid overfitting, just set the number of epochs as 2. Use a batch size of 10.

```
In [ ]: # resource: https://www.jianshu.com/p/d121ae396130?utm\_campaign=maleskine&utm\_content=note&utm\_medium=seo\_notes&utm\_source=wechat\_channel

model = Sequential()
model.add(Dense(50, activation = 'relu', input_dim = 23584))
model.add(Dropout(0.2))
model.add(Dense(50, activation = 'relu'))
model.add(Dropout(0.5))
model.add(Dense(50, activation = 'relu'))
model.add(Dropout(0.5))
model.add(Dense(1, activation = 'sigmoid'))

model.compile(loss = 'binary_crossentropy', optimizer = 'rmsprop', metrics = ['accuracy'])
model.fit(X_tr, Y_tr, epochs = 2, batch_size = 10)

Epoch 1/2
140/140 [=====] - 2s 9ms/step - loss: 0.7053 - accuracy: 0.4943
Epoch 2/2
140/140 [=====] - 1s 9ms/step - loss: 0.6841 - accuracy: 0.5241

ut[20]: <tensorflow.python.keras.callbacks.History at 0x1618ebd00>
```

ii. Report the train and test accuracies of this model

```
In [ ]: test_score = model.evaluate(X_t, Y_t, batch_size = 10)
print("test_score: ", mean(test_score))
train_score = model.evaluate(X_tr, Y_tr, batch_size = 10)
print("test_score: ", mean(train_score))

60/60 [=====] - 0s 1ms/step - loss: 0.6883 - accuracy: 0.5533
test_score: 0.6208227872848511
140/140 [=====] - 0s 1ms/step - loss: 0.6269 - accuracy: 0.6829
test_score: 0.6548685431480408
```


1.10 Developing a text classification model based on CNN + LSTM in Keras.

1.10.1 CNN based Text Classification:

One-Dimensional Convolutional Neural Network:

Although CNNs are mainly used for image data, they can also be applied to text data, as text also has adjacency information. Keras supports one-dimensional convolutions and pooling by the Conv1D and MaxPooling1D classes respectively.

i. After the embedding layer, insert a Conv1D layer. This convolutional layer has 32 feature maps, and each of the 32 kernels has size 3, i.e. reads embedded word representations 3 vector elements of the word embedding at a time. The convolutional layer is followed by a 1D max pooling layer with a length and stride of 2 that halves the size of the feature maps from the convolutional layer. The rest of the network is the same as the neural network above.

```
In [ ]: model = Sequential()
model.add(Embedding(5000, 32, input_length = L))
model.add(Conv1D(32,3))
model.add(MaxPooling1D(pool_size=2, strides=2))
model.add(Flatten())
model.compile('rmsprop', 'mse')
print(model.summary)

embed_text = model.predict(input_text)
print(model.layers[0].get_weights())
print(embed_text)

<bound method Model.summary of <tensorflow.python.keras.engine.sequential.Sequential object at 0x17c918a00>>
[array([[ -0.03625785, -0.00394065, -0.01474299, ..., 0.03426755,
         0.03313489, -0.03360778],
        [ 0.0382765 , 0.00851555, -0.02150942, ..., 0.0178039 ,
         0.01025463, -0.01649923],
        [ 0.00754521, 0.04487438, 0.00597467, ..., -0.03997531,
         -0.01019354, 0.04177039],
        ...,
        [-0.048969 , 0.03454694, -0.00401114, ..., 0.02102664,
         0.00111048, 0.02649759],
        [-0.03555398, -0.04463471, -0.0326045 , ..., -0.02369057,
         -0.0470906 , -0.02861959],
        [-0.03839446, 0.02004881, 0.02382291, ..., -0.01642703,
         -0.0390697 , -0.02226261]], dtype=float32)]
[[ 1.16310893e-02 1.59416627e-02 4.05392908e-02 ... 2.37603653e-02
  2.84849163e-02 -9.02268756e-03]
 [ 1.15169855e-02 4.81385097e-04 8.76351260e-03 ... 1.72630902e-02
  2.32951529e-02 9.26091988e-03]
 [-3.18915071e-03 5.93612343e-02 2.09271610e-02 ... 4.64434028e-02
 -8.77431966e-03 4.62909080e-02]
 ...,
 [ 1.15169855e-02 4.81385097e-04 8.76351260e-03 ... -3.15948762e-03
 -2.63488218e-02 2.62491424e-02]
 [ 1.15169855e-02 4.81385097e-04 8.76351260e-03 ... 3.67776155e-02
 2.89064739e-02 -2.69080908e-03]
 [ 1.15169855e-02 4.81385097e-04 8.76351260e-03 ... 3.42333689e-03
 4.80349222e-03 -6.62009697e-05]]

In [ ]: Data1 = data[["num", "class"]]
Data1["emb_text"] = list(embed_text)
Data1["class"] = Data1["class"].replace(-1, 0)

<ipython-input-23-2e2d5b954f82>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
Data1["emb_text"] = list(embed_text)
<ipython-input-23-2e2d5b954f82>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
Data1["class"] = Data1["class"].replace(-1, 0)

In [ ]: Train = Data1[Data1["num"] < 700]
Test = Data1[Data1["num"] >= 700]
X_tr = np.array(list(Train["emb_text"]))
X_t = np.array(list(Test["emb_text"]))
Y_tr = Train["class"]
Y_t = Test["class"]

In [ ]: model = Sequential()

model.add(Dense(50, activation = 'relu', input_dim = 11744))
model.add(Dropout(0.2))
model.add(Dense(50, activation = 'relu'))
model.add(Dropout(0.5))
model.add(Dense(50, activation = 'relu'))
model.add(Dropout(0.5))
model.add(Dense(1, activation = 'sigmoid'))

model.compile(loss = 'binary_crossentropy', optimizer = 'rmsprop', metrics = ['accuracy'])
model.fit(X_tr, Y_tr, epochs = 2, batch_size = 10)

Epoch 1/2
140/140 [=====] - 1s 6ms/step - loss: 0.7036 - accuracy: 0.5212
Epoch 2/2
140/140 [=====] - 1s 6ms/step - loss: 0.6949 - accuracy: 0.5070

Out[25]: <tensorflow.python.keras.callbacks.History at 0x17c647340>
```

1.10.2 LSTM based Text Classification:

1.10.3 Long Short-Term Memory Recurrent Neural Network

The structure of the LSTM we are going to use is shown in the following figure.

i. Each word is represented to LSTM as a vector of 32 elements and the LSTM is followed by a dense layer of 256 ReLUs. Use a dropout rate of 0.2 for both LSTM and the dense layer. Train the model using 10-50 epochs and batch size of 10.

```
In [ ]: model = Sequential()
model.add(Embedding(5000, 32, input_length = L))
model.add(LSTM(32))
model.add(Dropout(0.2))
model.add(Dense(256, activation = 'relu'))
model.add(Dropout(0.2))
model.add(Dense(1, activation = 'sigmoid'))

model.compile(loss = 'binary_crossentropy', optimizer = 'rmsprop', metrics = ['accuracy'])
#model.fit(X_tr, Y_tr, epochs = 10, batch_size = 10)
```

Type Markdown and LaTeX: α^2

```
In [ ]: test_score = model.evaluate(X_t, Y_t, batch_size = 10)
print("test_score: ", mean(test_score))
train_score = model.evaluate(X_tr, Y_tr, batch_size = 10)
print("test_score: ", mean(train_score))

WARNING:tensorflow:Model was constructed with shape (None, 737) for input KerasTensor(type_spec=TensorSpec(shape=(None, 737), dtype=tf.float32, name='embedding_4_input'), name='embedding_4_input', description="created by layer 'embedding_4_input'"), but it was called on an input with incompatible shape (10, 11744).
WARNING:tensorflow:Model was constructed with shape (None, 737) for input KerasTensor(type_spec=TensorSpec(shape=(None, 737), dtype=tf.float32, name='embedding_4_input'), name='embedding_4_input', description="created by layer 'embedding_4_input'"), but it was called on an input with incompatible shape (10, 11744).
60/60 [=====] - 21s 333ms/step - loss: 0.6956 - accuracy: 0.1632
test_score: 0.5965770184993744
140/140 [=====] - 48s 344ms/step - loss: 0.6932 - accuracy: 0.5000
test_score: 0.5965765416622162
```