```
In [407]: def plot_confusion_matrix(cm, classes, cmap=plt.cm.Blues):
       cm = cm.astype("float") / cm.sum(axis=1)
       plt.figure(figsize=(10,10))
       plt.imshow(cm, interpolation="nearest", cmap=cmap)
       plt.colorbar()
       tick_marks = np.arange(len(classes))
       plt.xticks(tick_marks, classes, fontsize=20, rotation=45)
      plt.yticks(tick_marks, classes, fontsize=20)
       for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
        color = "white" if cm[i,j] > 0.5 else "black"
        plt.text(j, i, format(cm[i, j], ".2f"), horizontalalignment="center", verticalalignm
     ent="center", color=color, fontsize=40)
      plt.ylabel("Correction", fontsize=15)
      plt.xlabel("Prédiction", fontsize=15)
In [408]: # nettoyage des données (ponctuations)
     stopWords = set(stopwords.words("english"))
     def cleanText(text):
      forbidden = {",","@",";","/","-",":",".","!","?", "#","\"","(",")","\'","/","/","-",
     "...", "&"}
       res = str(text)
      if res is None:
        return ""
      for elm in forbidden:
        res = res.replace(elm, "")
      if len(res.split()) >= 30:
        res = " ".join(res.split()[0:30])
      for elm in forbidden:
        res = res.replace(" ", " ")
       return res
     Traitement des données
     Dataset 1
In [409]: # chargement des deux tables
     real_news = pd.read_csv("./dataset1/True.csv")
     fake_news = pd.read_csv("./dataset1/Fake.csv")
     # vrai = 1, faux = 0
     real_news["label"] = 1
     fake_news["label"] = 0
     # création du dataset complet
     dataframe1 = pd.concat([real_news, fake_news])
     print(f"Nombre de références : {dataframe1.title.count()}")
     print(f"Nombre de fake news : {fake_news.title.count()}")
     print(f"Nombre de vraies news : {real_news.title.count()}")
     # ici on ne s'intéresse qu'au titre et au label
     del dataframe1["title"]
     del dataframe1["subject"]
     del dataframe1["date"]
     dataframe1["text"] = dataframe1["text"].apply(cleanText)
     dataframe1.sample(10)
     Nombre de références : 44898
     Nombre de fake news : 23481
     Nombre de vraies news : 21417
Out[409]:
                        text label
     13047
          I know there are only 60 days left to make our...
     7876
                           0
          Marijuana legalization is on the cusp in Ameri...
     2739
         WASHINGTON Reuters With talks to renegotiate t...
        WASHINGTON Reuters President Donald Trumps cal...
     1860
                           1
     18764
        BARCELONA Reuters Catalan separatists urged su...
                           1
       WASHINGTON Reuters US President Donald Trump s...
     1842
                           1
     16893
          The Transportation Security Administration TSA...
     14516
        JOHANNESBURG Reuters South Africa s Justice Mi...
                           1
     9780
         During game 1 of the WNBA Finals the LA Sparks...
     3038
        COLORADO SPRINGS Colo Reuters Officials with t...
     Dataset 2
     dataframe2 = pd.read_csv("./dataset2/train.csv")
     del dataframe2["id"]
     del dataframe2["author"]
     del dataframe2["title"]
     dataframe2["text"] = dataframe2["text"].apply(cleanText)
     dataframe2.sample(10)
Out[410]:
                        text label
        WASHINGTON — Attorney General Jeff Sessions fa...
                           0
     5801
        20 Views November 03 2016 GOLD KWN King World ...
     16070
     11234
           Print Protesters blocked the upper level of th...
     2134
         « on Today at 083635 PM » AIG Quadruples Limit...
     4578
        Wed 26 Oct 2016 1945 UTC © The Canary Leaked d...
     15668
         Good morning Were trying something new this we...
     17032
          License DMCA With the US more interested in fi...
     3190 STAFFORD SPRINGS Conn — Sandra Miller was at w...
     4575
           Veteran political strategist Pat Caddell talke...
     10752
          Leaked Podestas Satanic Spirit Cooking Dinner ...
In [411]:
     # on sépare les données en données d'entraînement et données de test (80% et 20%)
     x_train, x_test, y_train, y_test = train_test_split(dataframe1["text"], dataframe1["label"],
     test_size=0.99, random_state = 42)
     print(f"Données d'entrainement : {len(x_train)}")
     print(f"Données de test : {len(x_test)}")
     Données d'entrainement : 448
     Données de test : 44450
In [412]: max_features = 10000 # taille max du vocab
     maxlen = 70 # taille max de séquence
In [413]: # vectorisation naïve en "one-hot"
     tokenizer = Tokenizer(num_words=max_features)
     tokenizer.fit_on_texts(pd.concat([dataframe1["text"], dataframe2["text"]]))
In [414]: # vectorisation des données d'entraînement
     x_train = tokenizer.texts_to_sequences(x_train)
     x_train = pad_sequences(x_train, maxlen=maxlen)
In [415]: # vectorisation des données de test
     x_test = tokenizer.texts_to_sequences(x_test)
     x_test = pad_sequences(x_test, maxlen=maxlen)
     Modèle Perceptron
     batch_size = 128 # la batch_size n'influe pas sur le résultat non plus
     nb_epochs = 100 # le nombre d'époques n'influe pas sur le résultat
     embedded_dim = 100 # idem
In [417]: model = Sequential()
     model.add(Embedding(max\_features, output\_dim=embedded\_dim, input\_length=maxlen, trainable=Tr
     model.add(Flatten())
     model.add(Dense(256, activation="relu"))
     model.add(Dropout(0.1))
     model.add(Dense(128, activation="relu"))
     model.add(Dense(1, activation="sigmoid"))
     model.compile(loss="binary_crossentropy", optimizer=Adam(lr=0.0002), metrics=["accuracy"])
     model.summary()
     Model: "sequential_16"
                               Param #
     Layer (type)
                   Output Shape
     embedding_16 (Embedding)
                               1000000
                   (None, 70, 100)
     flatten_14 (Flatten)
                   (None, 7000)
                               0
     dense_25 (Dense)
                               1792256
                   (None, 256)
     dropout_15 (Dropout)
                               0
                   (None, 256)
     dense_26 (Dense)
                   (None, 128)
                               32896
     dense_27 (Dense)
                               129
                   (None, 1)
     ______
     Total params: 2,825,281
     Trainable params: 2,825,281
     Non-trainable params: 0
In [418]: progress = model.fit(x_train, y_train, batch_size=batch_size, epochs=nb_epochs)
     C:\Users\Shadow\AppData\Roaming\Python\Python37\site-packages\tensorflow\python\framework\ind
     exed_slices.py:434: UserWarning: Converting sparse IndexedSlices to a dense Tensor of unknown
     shape. This may consume a large amount of memory.
     "Converting sparse IndexedSlices to a dense Tensor of unknown shape."
     Epoch 2/100
     Epoch 3/100
     Epoch 4/100
     Epoch 5/100
     Epoch 6/100
     Epoch 7/100
     Epoch 8/100
     Epoch 9/100
     Epoch 10/100
     Epoch 11/100
     Epoch 12/100
     Epoch 13/100
     Epoch 14/100
     Epoch 15/100
     Epoch 16/100
     Epoch 17/100
     Epoch 18/100
     Epoch 19/100
     Epoch 20/100
     Epoch 21/100
     Epoch 22/100
     Epoch 23/100
     Epoch 24/100
     Epoch 25/100
     Epoch 26/100
     Epoch 27/100
     Epoch 28/100
     Epoch 29/100
     Epoch 30/100
     Epoch 31/100
     Epoch 32/100
     Epoch 33/100
     Epoch 34/100
     Epoch 35/100
     Epoch 36/100
     Epoch 37/100
     Epoch 38/100
     Epoch 39/100
     Epoch 40/100
     Epoch 41/100
     Epoch 42/100
     Epoch 43/100
     Epoch 44/100
     Epoch 45/100
     Epoch 46/100
     Epoch 47/100
     Epoch 48/100
     Epoch 49/100
     Epoch 50/100
     Epoch 51/100
     Epoch 52/100
     Epoch 53/100
     Epoch 54/100
     Epoch 55/100
     Epoch 56/100
     Epoch 57/100
     Epoch 58/100
     Epoch 59/100
     Epoch 60/100
     Epoch 61/100
     Epoch 62/100
     Epoch 63/100
     Epoch 64/100
     Epoch 65/100
     Epoch 66/100
     Epoch 67/100
     Epoch 68/100
     Epoch 69/100
     Epoch 70/100
     Epoch 71/100
     Epoch 72/100
     Epoch 73/100
     Epoch 74/100
     Epoch 75/100
     Epoch 76/100
     Epoch 77/100
     Epoch 78/100
     Epoch 79/100
     Epoch 80/100
     Epoch 81/100
     Epoch 82/100
     Epoch 83/100
     Epoch 84/100
     Epoch 85/100
     Epoch 86/100
     Epoch 87/100
     Epoch 88/100
     Epoch 89/100
     Epoch 90/100
     Epoch 91/100
     Epoch 92/100
     Epoch 93/100
     Epoch 94/100
     Epoch 95/100
     Epoch 96/100
     Epoch 97/100
     Epoch 98/100
     Epoch 99/100
     448/448 [=====
                :=========] - Os 63us/step - loss: 2.3656e-04 - accuracy: 1.0000
     Epoch 100/100
     Analyse des résultats
In [419]: # art plastique du turfu featuring le poto matplotlib
     from matplotlib import pyplot as plt
In [420]: figure, ax = plt.subplots(1, 2)
     figure.set_size_inches(20,10)
     ax[0].plot(progress.history["loss"])
     ax[1].plot(progress.history["accuracy"])
     plt.show()
     0.7
                             1.0
     0.6
                             0.9
     0.5
                             0.8
     0.4
     0.3
     0.2
                             0.6
     Ci dessus : mise en évidence que le nombre d'époques influe peu : le modèle devient très précis sur le dataset au bout de 20
     époques
    dataframe = pd.concat([dataframe1, dataframe2])
In [421]:
    to_predict = dataframe2.sample(10000)
In [434]:
In [435]: X_test = to_predict["text"]
In [436]: Y_test = to_predict["label"]
In [437]: X_test = tokenizer.texts_to_sequences(X_test)
     X_test = pad_sequences(X_test, maxlen=maxlen)
     predictions = model.predict(X_test)
In [438]:
     threshold = 0.5 #la valeur du threshold n'a pas d'importance, les réseaux perceptrons ne fon
     ctionnent pas sur ce problème
     predictions[predictions > threshold] = 1
     predictions[predictions <= threshold] = 0</pre>
In [439]: cm_plot_labels = ["Vraie Info", "Fausse Info"]
     cm = confusion_matrix(Y_test, predictions)
     plot_confusion_matrix(cm, cm_plot_labels)
                                          0.8
                 0.81
                                          0.7
                              0.19
       Vraie Info
                                          0.6
     Correction
                                          0.5
                                          0.4
                              0.08
                 0.93
      Fausse Info
                                          0.3
                                          0.2
                                          0.1
                        Prédiction
     Cela démontre bien qu'un réseau perceptron n'est pas capable de généraliser la notion de fake news sur un dataset sur
     lequel il n'a pas été entraîné.
```

Inefficacité du Perceptron

from nltk.corpus import stopwords

from keras.models import Sequential

from keras.optimizers import Adam

from sklearn.model_selection import train_test_split

from keras.preprocessing.sequence import pad_sequences

from keras.layers import Conv2D, Dropout, Dense, Embedding, Flatten, Reshape

from sklearn.metrics import confusion_matrix
from keras.preprocessing.text import Tokenizer

In [406]: import pandas as pd

import numpy as np
import itertools

Boite à outils