

# PROJET 1

**INTRODUCTION AU METIER D'INGENIEUR IA:**  
MISE EN OEUVRE DE LA FONCTION POLYGLOTTE  
D'UN CHATBOT VIA DES SERVICES AZURE

**#AZURE COGNITIVE  
SERVICES #PYTHON**

## Ingénieur IA

Développez et intégrez des algorithmes de Deep Learning au sein d'un produit IA

Ouddane Nabil



# SOMMAIRE

---

## Projet 1

Découvrez le métier d'Ingénieur IA

### A. INTRODUCTION

1. Contexte
2. Objectifs

### B. PREREQUIS AU PROJET

1. Création d'un compte azure
2. Python
3. Création de variables d'environnement permettant l'accès au coffre fort via l'application P1APP autorisée
4. Insertion des clés API des services dans le coffre fort AZURE KEYVAULT

### C. SCRIPT

1. Fonctionnement du script
2. Import des données linguistiques wikipedia
3. Extraction des paragraphes des 5 langues les plus parlées
4. EXTRACTION DE N PARAGRAPHES ALEATOIRES POUR NE PAS EXPLOSER LES COMPTEURS AZURE
5. A) Authentification sécurisée au service text analytics & B) Authentification sécurisée au service translator
6. A) Détection du langage via l'IA TEXT ANALYTICS d'AZURE & B) Détection du langage via l'IA TRANSLATOR/DETECT

### D. CONCLUSION

1. Les Conclusions

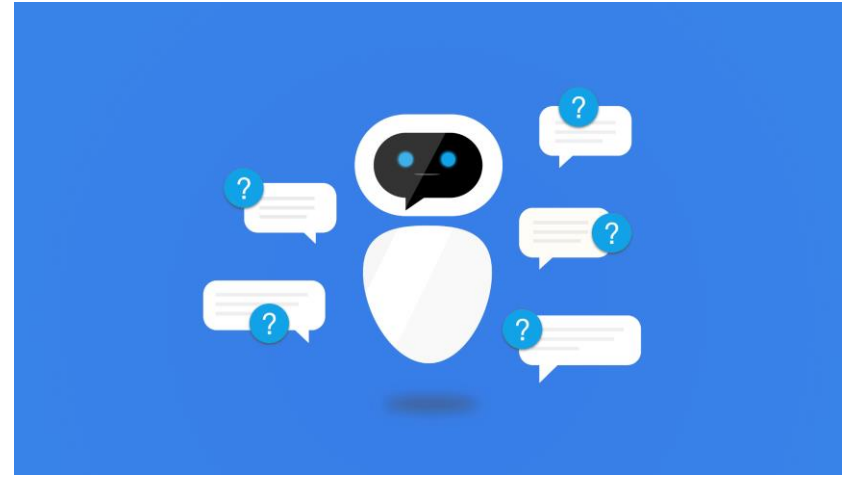


# INTRODUCTION

# 1. Contexte

---

- **ENJEU global:** Conception d'un chatbot polyglotte, fiable et pertinent pour fluidifier l'expérience client du site web de la banque CREDIT
- **ENJEU DU P1:** Mise en Oeuvre de la fonction polyglotte du chatbot
  - Détection du langage par les modèles pré-entraînés des fonctions cognitives du CLOUD Microsoft AZURE
- **SUGGESTION DE RESSOURCES:**
  - Jeu de données test [Wikipedia Language Identification Database](#)
  - API CURL de la ressource TRANSLATOR text: <https://docs.microsoft.com/en-us/azure/cognitive-services/translator/reference/v3-0-detect>



## 2. Objectifs

---

- **SCRIPT**: délivrer du code présentable permettant de
  - Modifier facilement l'input à tester
  - Se connecter de façon sécurisée au service Azure:
    - Les clés de connexion ne doivent pas être visible
  - Récupérer la prédiction du modèle de détection de langue

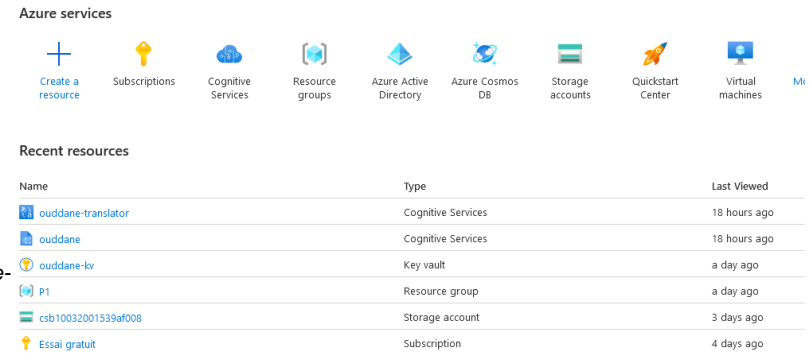




# PREREQUIS AU PROJET

# 1. Création d'un compte azure

- Création de la ressource **Microsoft.CognitiveServicesTextAnalytics** dans azure
  - Ce Service possède une fonction de détection de langage bien que non suggérée ainsi qu'une librairie python
- Création de la ressource **Microsoft.CognitiveServicesTranslator** dans azure
  - Ce Service possède une fonction de détection de langage suggérée
  - Son utilisation se fait par une API CURL
- Création de la ressource **Microsoft azure key vault** dans azure
  - Permet de stocker de manière sécurisée les clés d'accès aux API des services : <https://docs.microsoft.com/en-us/azure/key-vault/secrets/quick-create-python>
  - Utilisation du bash azure CLI en ligne pour la création du coffre fort: <https://shell.azure.com/>
    - Création du coffre fort en ligne de commande:  
`az group create --name KeyVault-PythonQS-rg --location westeurope`  
`az keyvault create --name <your-unique-keyvault-name> ouddane-kv --resource-group KeyVault-PythonQS-rg`
    - Création d'une application autorisée **P1app** pour le coffre fort dans active directory et paramétrage des accès au coffre fort: <https://www.c-sharpcorner.com/article/how-to-access-azure-key-vault-secrets-through-rest-api-using-postman/>



## 2. Python

---

- **INSTALLATION DE LA SUITE ANACONDA**

- <https://www.anaconda.com/>

- **CHOIX DE L'ENVIRONNEMENT JUPYTER NOTEBOOK**

- Adapté pour les data scientist

- **INSTALLATION DE LIBRAIRIES AZURE VIA**

**JUPYTER:**<https://jakevdp.github.io/blog/2017/12/05/installing-python-packages-from-jupyter/>

- azure-ai-textanalytics
    - Service azure cognitive services permettant la detection de langue
  - azure-identity & azure-keyvault-secrets
    - Service permettant de gérer l'accès au coffre fort azure et la connexion sécurisée aux services azure

### STEP 0.1: INSTALLATION DES LIBRAIRIES AZURE NECESSAIRES VIA JUPYTER

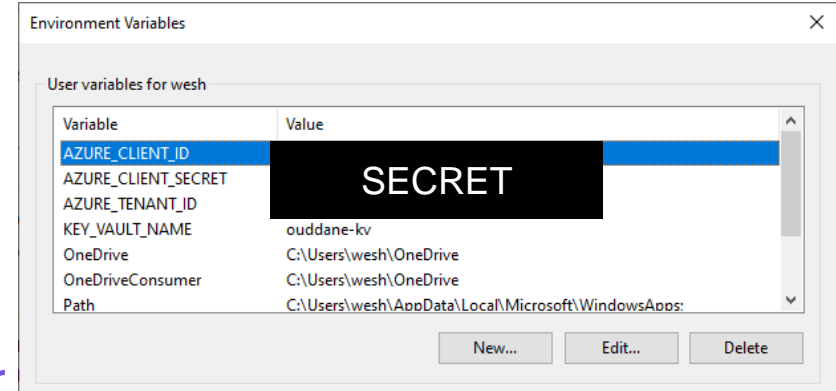
Installation de **azure-ai-textanalytics** pour l'analyse de texte IA Installation de **azure-identity** et **azure-keyvault-secrets** pour les accès au coffre fort AZURE <https://jakevdp.github.io/blog/2017/12/05/installing-python-packages-from-jupyter/>

```
: 1 #aminata1902@hotmail.fr
2 #STEP 0.1
3 #prerequisite (one time)
4
5 import sys
6 !{sys.executable} -m pip install azure-ai-textanalytics
7
8 #creation de la ressource Microsoft.CognitiveServicesTextAnalytics dans azure
9 #https://docs.microsoft.com/en-us/azure/cognitive-services/text-analytics/quickstarts/client-libraries-rest-api?tabs=ver
10 #https://azure.microsoft.com/fr-fr/blog/getting-started-with-cognitive-services-language-understanding-container/
11
12 !{sys.executable} -m pip install azure-identity
13 !{sys.executable} -m pip install azure-keyvault-secrets
14
```



### 3. Création de variables d'environnement permettant l'accès au coffre fort via l'application P1APP autorisée

- Création des clés suivantes disponibles sur son compte azure paramétré
  - AZURE\_CLIENT\_ID
  - AZURE\_CLIENT\_SECRET
  - AZURE\_TENANT\_ID
  - KEY\_VAULT\_NAME
- Insertion ou modification de ses valeurs par un code python si besoin



- la fonction **os.environ** map les variables environnementales une fois pour toutes à l'import de la librairie os. Les changements de l'environnement opérés après ce mapping ne sont pas répercutés dans os.environ, à part les modifications directes de os.environ dans le code.

```
1 #STEP 0.3
2 #environ var set up
3
4 import os
5
6 #something weird
7 #need first to set up the environment variables on the system (control panel)
8
9
10 os.environ['KEY_VAULT_NAME'] = 'ouddane-kv'
11 os.environ['AZURE_TENANT_ID'] = 'SECRET'
12 #https://docs.microsoft.com/fr-fr/active-directory-how-to-find-tenant
13 os.environ['AZURE_CLIENT_ID'] = 'SECRET'
14 #https://www.c-sharpcorner.com/active-directory-how-to-find-tenant
15 os.environ['AZURE_CLIENT_SECRET'] = 'SECRET'
16
17 #import os
18 #print(os.getenv('KEY_VAULT_NAME'))
19 #print(os.getenv('AZURE_TENANT_ID'))
20 #print(os.getenv('AZURE_CLIENT_ID'))
21 #print(os.getenv('AZURE_CLIENT_SECRET'))
22 #print(os.environ)
```

## 4. Insertion des clés API des services dans le coffre fort AZURE KEYVAULT

---

- Création d'un script python permettant de demander les clés puis de les stocker
  - L'accès est complètement sécurisé
  - **STEP 0.4: STOCKAGE DE LA KEY cognitives services DANS LE COFFRE FORT AZURE KEYVAULT**

<https://docs.microsoft.com/fr-fr/azure/key-vault/secrets/quick-create-python>

<https://www.c-sharpcorner.com/article/how-to-access-azure-key-vault-secrets-through-rest-api-using-postman/>

```
: 1 #STEP 0.4
2 #code to store keys within azure keyvaults
3 import os
4 import cmd
5 from azure.keyvault.secrets import SecretClient
6 from azure.identity import DefaultAzureCredential
7
8 keyVaultName = os.getenv('KEY_VAULT_NAME')
9 KVUri = f"https://{keyVaultName}.vault.azure.net"
10
11 #https://docs.microsoft.com/en-us/python/api/overview/azure/identity-readme?view=azure-python
12
13 credential = DefaultAzureCredential()
14 client = SecretClient(vault_url=KVUri, credential=credential)
15
16 secretName = input("Input a name for your secret > ")
17 secretValue = input("Input a value for your secret > ")
18
19 print(f"Creating a secret in {keyVaultName} called '{secretName}' with the value '{secretValue}' ...")
20
21 client.set_secret(secretName, secretValue)
22
23 print(" done.")
```

Input a name for your secret > test

Input a value for your secret > test

Creating a secret in ouddane-kv called 'test' with the value 'test' ...

done.

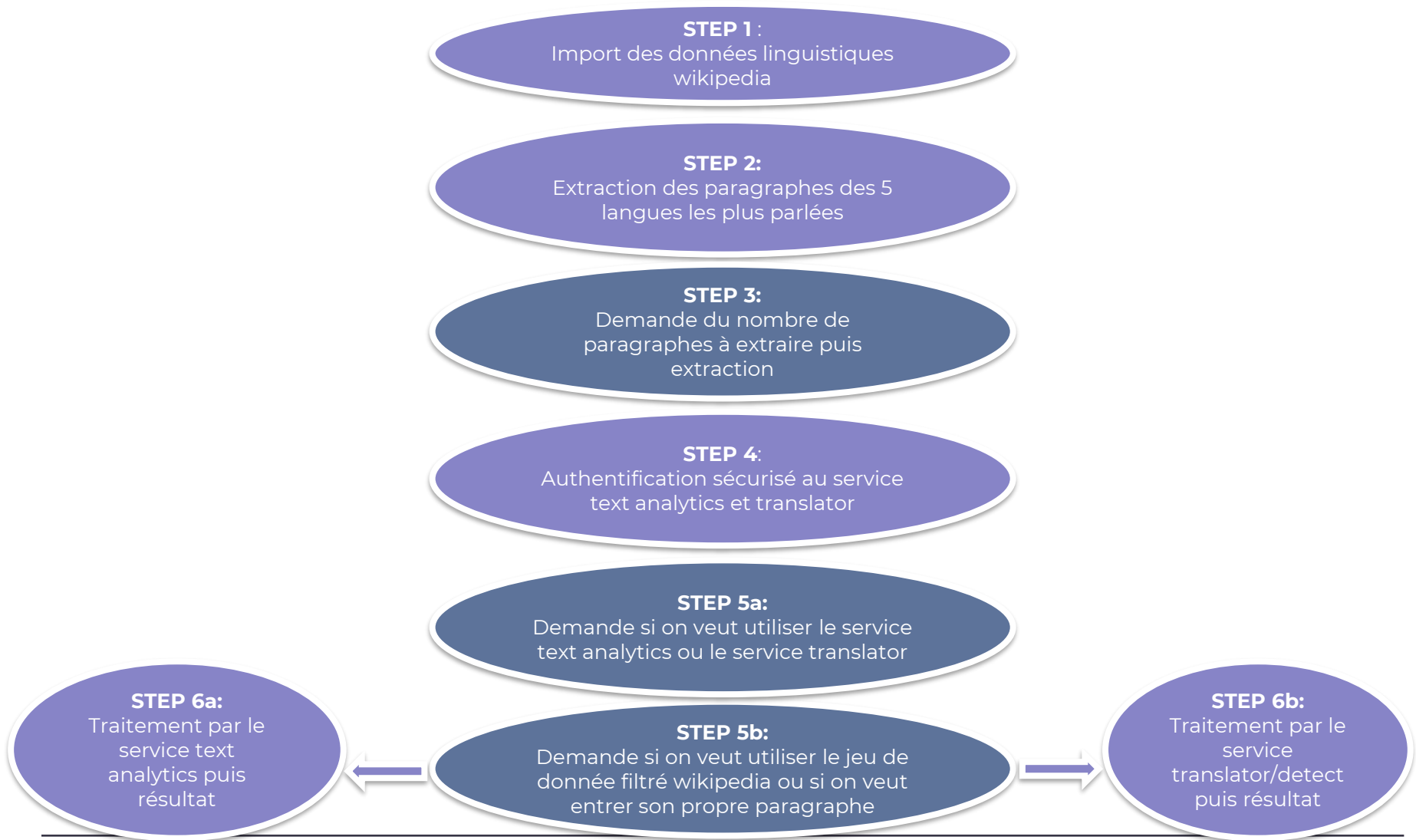
---



# 1. Fonctionnement du script

## Possibilité d'utiliser 2 services azure différents

---



## 2. Import des données linguistiques wikipedia

---

- Utilisation des fonctions de la librairie pandas
  - ne pas oublier la petite option quoting=3 dans read.table() qui fait la différence pour télécharger tous les lignes proprement

```
1 #step 1
2 #wikipedia data languages reading
3 #test avec les dataframes de pandas
4
5 #When reading the txt file with pandas' read_csv, some line loss occurred.
6 #The reason is that there is a single English quotation mark in a line, causing the \n line break to be invalid.
7 #Multiple lines are concatenated until the next single quote is encountered
8 #add option quoting=3 to tackle the issue :QUOTE_NONE
9
10 import pandas as pd
11 import numpy as np
12
13 dflabel = pd.read_csv('input/labels.csv', delimiter = ';')
14 dfx1 = pd.read_table('input/x_train.txt', delimiter = None, header=None, encoding='utf-8', quoting=3)
15 dfy1 = pd.read_table('input/y_train.txt', delimiter = None, header=None, encoding='utf-8', quoting=3)
16 dfx2 = pd.read_table('input/x_test.txt', delimiter = None, header=None, encoding='utf-8', quoting=3)
17 dfy2 = pd.read_table('input/y_test.txt', delimiter = None, header=None, encoding='utf-8', quoting=3)
18
19 #row number quick checking
20 print(len(dfx1.index))
21 print(len(dfx2.index))
22 print(len(dfy1.index))
23 print(len(dfy2.index))
24 print(len(dflabel.index))
```

```
117500
117500
117500
117500
235
```

### 3. Extraction des paragraphes des 5 langues les plus parlées

---

- Utilisation des fonctions de la librairie pandas
  - Selon WIKIPEDIA les 5 langues les plus parlées sont les suivantes: english / mandarin / hindi / spanish / standard arabic: [https://en.wikipedia.org/wiki/List\\_of\\_languages\\_by\\_total\\_number\\_of\\_speakers](https://en.wikipedia.org/wiki/List_of_languages_by_total_number_of_speakers)

```
: 1 #step 2
2 #5 most spoken languages selection
3 #https://en.wikipedia.org/wiki/List_of_languages_by_total_number_of_speakers
4 #english / mandarin / hindi / spanish / standard arabic
5
6 #sentences and language merging
7 df1 = pd.concat([dfx1, dfy1], axis=1)
8 df1.columns=["para", "langcod"]
9 df2 = pd.concat([dfy2, dfy2], axis=1)
10 df2.columns=["para", "langcod"]
11
12 #5 languages extraction from labels
13 dflabel5mml=dflabel.loc[(dflabel['English']=='Arabic') |
14                        (dflabel['English']=='English') |
15                        (dflabel['English']=='Hindi') |
16                        (dflabel['English']=='Spanish') |
17                        (dflabel['English']=='Standard Chinese')]
18 dflabel5mml2=dflabel5mml.rename(columns={'Label':'langcod'})
19 dflabel5mml3=dflabel5mml2[['langcod', 'English']]
20
21 #sentences selection by most spoken languages-> merging
22 df1s = pd.merge(df1, dflabel5mml3, how="inner", on=["langcod"])
23 df2s = pd.merge(df2, dflabel5mml3, how="inner", on=["langcod"])
24
25 #row number quick checking
26 print(len(df1s.index))
27 print(len(df2s.index))
```

2500

2500

---

## 4. EXTRACTION DE N PARAGRAPHES ALEATOIRES POUR NE PAS EXPLOSER LES COMPTEURS AZURE

- Demande du nombre de paragraphes à extraire des données sources préfiltrées
  - Nombre minimum de 5 pour obtenir les 5 langues.
  - Tirage aléatoire de N paragraphes provenant des 2500 paragraphes préfiltrés.
    - Tirage équitable dans les 5 langues
  - Récupération également des langue correspondantes d'après wikipedia

```
1 #step 3
2 #extraction of n random sentences form dataframe
3 import random
4
5 def rand_shrink_df(df):
6
7     size = input("Input a number of paragraphs, bigger than 5, that you want to test ,\
8 keep in mind that azure is not free > ")
9     sizeint = int(size)
10    #Generate random numbers
11    #randomlist = random.sample(range(0, len(df.index)-1), sizeint)
12    randomlist = random.sample(range(0, 500), sizeint)
13    #print(randomlist)
14    #trick to get all the language
15    factor = [(i % 5)*500 for i in range(sizeint)]
16    #print(factor)
17    #operation of 2 list with zip
18    product = [x+y for x,y in zip(randomlist,factor)]
19    #print(product)
20
21
22
23    for i in range(sizeint):
24        #attention pour que le dataframe reste en row, penser aux double crochets
25        dfa=df.iloc[[product[i]],:]
26        #print(i)
27        if i==0 :
28            dfaa=dfa
29            #print(dfaa)
30            #print(len(dfaa.index))
31        else:
32            dfaa=dfa.append(dfa)
33            #print(len(dfa.index))
34    return dfaa
35
36 df2az=rand_shrink_df(df1s)
37 #print(len(df2az.index))
38 print(df2az)
39
40 #faire en sorte que toutes les langues apparaissent dans le tirage
```

Input a number of paragraphs, bigger than 5, that you want to test ,keep in mind that azure is not free >

## 5. A) Authentification sécurisée au service text analytics

---

- Connexion sécurisé au service azure keyvault
  - Récupération sécurisé de la Key API permettant la connexion au service text analytics
  - Aucune clé ne figure dans le code, des variables d'environnement sont utilisées
- Connexion sécurisée au service text analytics grâce a la key API

```
1  #step 4a
2  #authenticate the client
3  #azure key vault (coffre fort) to securely store keys
4  from azure.keyvault.secrets import SecretClient
5  from azure.identity import DefaultAzureCredential
6
7  keyVaultName = os.getenv('KEY_VAULT_NAME')
8  KVUri = f"https://{keyVaultName}.vault.azure.net"
9
10 #https://docs.microsoft.com/en-us/python/api/overview/azure/identity-readme?view=azure-python
11
12 credential = DefaultAzureCredential()
13 #get the needed environmental variables to connect to the keyvault
14 client = SecretClient(vault_url=KVUri, credential=credential)
15
16 secretName = "keytextanalytics"
17 endpoint = "https://cuddane.cognitiveservices.azure.com/"
18
19 #https://docs.microsoft.com/en-us/azure/cognitive-services/text-analytics/quickstarts/client-libraries-rest-api?tabs=ver
20 # use this code if you're using SDK version is 5.0.0
21 from azure.ai.textanalytics import TextAnalyticsClient
22 from azure.core.credentials import AzureKeyCredential
23
24 def authenticate_client():
25     key=client.get_secret(secretName)
26     ta_credential = AzureKeyCredential(key.value)
27     text_analytics_client = TextAnalyticsClient(
28         endpoint=endpoint,
29         credential=ta_credential)
30     return text_analytics_client
31
32 client = authenticate_client()
```



## 5. B) Authentification sécurisée au service translator

- Connexion sécurisé au service azure keyvault
  - Récupération sécurisé de la Key API permettant la connexion au service translator
  - Aucune clé ne figure dans le code, des variables d'environnement sont utilisées
- La connexion sécurisée au service translator se fait en envoyant la requête POST de détection de langage en passant la clé API comme paramètre headers

```
1 #step 6
2 #language detection by translator
3 import os
4 import cmd
5 import requests, uuid, json
6
7 #authenticate the client
8 #azure key vault (coffre fort) to securely store keys
9 from azure.keyvault.secrets import SecretClient
10 from azure.identity import DefaultAzureCredential
11
12 def azure_translator_language_detection(df):
13
14
15     #https://docs.microsoft.com/en-us/python/api/overview/azure/identity-readme?view=azure-python
16     credential = DefaultAzureCredential()
17     #get the needed environmental variables to connect to the keyvault
18     keyVaultName = os.getenv('KEY_VAULT_NAME')
19     KVUri = f"https://{keyVaultName}.vault.azure.net"
20     client = SecretClient(vault_url=KVUri, credential=credential)
21
22     #get the key from the azure key vault to connect to translator
23     secretName = "keytranslator"
24     key=client.get_secret(secretName)
25
26     #https://docs.microsoft.com/en-us/azure/cognitive-services/translator/quickstart-translator?ta
27
28     # Add your subscription key and endpoint
29     subscription_key = key.value
30     endpoint = "https://api.cognitive.microsofttranslator.com/detect?api-version=3.0"
31     location = "francecentral"
32
33     headers = {
34         'Ocp-Apim-Subscription-Key': subscription_key,
35         'Content-type': 'application/json',
36         'Ocp-Apim-Subscription-Region': location,
37         'X-ClientTraceId': str(uuid.uuid4())
38     }
```

# 6. A) Détection du langage via l'IA TEXT ANALYTICS d'AZURE

- Demande du type d'input à tester
  - Paragraphe 'on the fly': On copie colle un paragraphe et on récupère la langue après requête
  - table préfiltrée des paragraphes wikipédia : On récupère un pourcentage moyen, très bon, de matching avec les valeurs de langue wikipédia et on affiche les paragraphes posant problème.

```
1 #step 5
2 #language detection by text analytics
3 def azure_text_analytics_language_detection(client,df):
4
5     own= input("Do you want to test your own sentence? write Y or N > ")
6     if own=="Y":
7         doc= [input("paste your sentence > ")]
8         response = client.detect_language(documents = doc, country_hint = "")[0]
9         print("According to AZURE text analytics, the language of your sentence is ",response.primary_language.name )
10    else:
11        df["azure"] = ""
12        df["flag"]= 0
13        for i in range(len(df.index)):
14            try:
15                documents = [df.iloc[i,0]]
16                response = client.detect_language(documents = documents, country_hint = "")[0]
17                df.iloc[i,3]=response.primary_language.name
18                if (df.iloc[i,3]== "Chinese_Traditional") | (df.iloc[i,3]== "Chinese_Simplified"):
19                    df.iloc[i,3]= "Standard Chinese"
20                if df.iloc[i,3]== df.iloc[i,2]:
21                    df.iloc[i,4]=1
22            else:
23                df.iloc[i,4]=0
24            except Exception as err:
25                print("Encountered exception. {}".format(err))
26
27        print("La detection fonctionne à",df["flag"].mean()*100,"%")
28        if df["flag"].mean()<1:
29            print("L'erreur provient peut etre de la donnée source")
30            df0=df.loc[(df['flag']==0)]
31            print(df0)
32
33    azure_text_analytics_language_detection(client,df2az)
```

Do you want to test your own sentence? write Y or N > Y  
paste your sentence > Au menu de la journée, deux huitièmes de finale : le premier entre les Pays-Bas et la République tchèque, le second entre la Belgique et le Portugal. Chaque matin, retrouvez l'actualité de la seizième édition de l'Euro de football.

According to AZURE, the language of your sentence is French

## 6. B) Détection du langage via l'IA TRANSLATOR/DETECT

- Demande du type d'input à tester
  - Paragraphe 'on the fly':  
On copie colle un paragraphe et on récupère la langue après requête. On peut facilement afficher le pourcentage de confiance de la détection
- table préfiltrée des paragraphes wikipedia :  
On récupère un pourcentage moyen, très bon, de matching avec les valeurs de langue wikipédia , un pourcentage moyen de confiance de la détection fourni par azure et on affiche les paragraphes posant problème .

```
40 own= input("Do you want to test your own sentence? write Y or N > ")
41 if own=="Y":
42     doc= input("paste your sentence > ")
43     # You can pass more than one object in json body.
44     body = [{
45         'text': doc
46     }]
47     ## lib requests
48     ## post sends a POST request (http) to the specified url and is used when you want to send some data.
49     request = requests.post(endpoint, headers=headers, json=body)
50     ##response.json() returns a JSON object of the result / Whenever we make a request to a specified URL
51     response = request.json()
52
53     #json.dumps() function converts a Python object into a json string
54     response2 = json.dumps(response)
55     #json.loads() method can be used to parse a valid JSON string and convert it into a Python Dictionary
56     dictresponse = json.loads(response2)[0]
57
58     #print(dictresponse)
59     print("According to AZURE text analytics, the language of your sentence is ", dictresponse['language'])
60
61 else:
62     #init flag
63     df["azure"] = ""
64     df["flag"] = 0
65     df["azscore"] = 0
66     for i in range(len(df.index)):
67         # try:
68         doc = df.iloc[i,0]
69         body = [{
70             'text': doc
71         }]
72         request = requests.post(endpoint, headers=headers, json=body)
73         response = request.json()
74         response2 = json.dumps(response)
75         dictresponse = json.loads(response2)[0]
76         #print(dictresponse)
77         df.iloc[i,3] = dictresponse['language']
78         df.iloc[i,5] = dictresponse['score']
79         if ((df.iloc[i,3]=="en") & (df.iloc[i,2]=="English"))|
80             ((df.iloc[i,3]=="ar") & (df.iloc[i,2]=="Arabic"))|
81             ((df.iloc[i,3]=="es") & (df.iloc[i,2]=="Spanish"))|
82             ((df.iloc[i,3]=="hi") & (df.iloc[i,2]=="Hindi"))|
83             ((df.iloc[i,3]=="zh-Hant") & (df.iloc[i,2]=="Standard Chinese"))|
84             ((df.iloc[i,3]=="zh-Hans") & (df.iloc[i,2]=="Standard Chinese"))):
85             df.iloc[i,4]=1
86         else:
87             df.iloc[i,4]=0
88         # except Exception as err:
89         #     print("Encountered exception. {}".format(err))
90
91     print("La detection fonctionne à",df["flag"].mean()*100,"% selon notre indicateur et le score moyen")
92     if df["flag"].mean()<1:
93         print("L'erreur provient peut etre de la donnée source: mauvais flag ou plusieurs langues dsur :
94         #65417 de train.txt: plusieurs langues
95         df0=df.loc[(df['flag']==0)]
96         print(df0)
97
98 azure_translator_language_detection(df2az)
```

# **D** CONCLUSION

# 1. Les Conclusions

---

- Fiabilité des services de détection de langue d'azure
    - Aussi bien pour le service text analytics que pour le service translator
  - Les données sources wikipedia ne sont pas 100% propres
    - Paragraphes avec plusieurs langues
    - Langue associée fausse : exemple: ligne 11251 de train.txt reference hindi alors que c est de l'anglais
  - Projet
    - Formateur sur python , azure et d'autres librairies
    - La réalisation répond à la demande et va au delà en testant 2 services différents
-