

The background of the slide features a complex network diagram. It consists of numerous nodes of varying sizes, some colored in dark blue, light blue, and grey, connected by a web of thin, light grey lines. Several nodes are highlighted with larger, concentric circles around them, suggesting a focus or a central point in the network. The overall aesthetic is modern and technological.

DÉTECTER LES BAD BUZZ* GRÂCE AU DEEP LEARNING

Voahangy Joan ALEONARD – 11/05/2021

**phénomène de bouche-à-oreille négatif sur le net*

AGENDA DU JOUR



PROJET ET
DÉMARCHE



PRÉ-TRAITEMENT
DES DONNÉES



PRÉSENTATION DES
APPROCHES



CHOIX MODÈLE ET
DÉPLOIEMENT



PROJET ET DÉMARCHE

CONTEXTE PROJET ET DÉMARCHE



Problématique

Surveiller la réputation sur les réseaux sociaux

Prédire le sentiment associé à un tweet



Livrer un prototype fonctionnel du modèle



Source de travail

Utiliser des données open source d'analyse de sentiment

Trouver des tweets annotés du sentiment exprimé (👍 ou 👎)



Jeu de données : [Sentiment 140](#)



Démarche

Analyser et prétraiter les données

Explorer les différentes approches

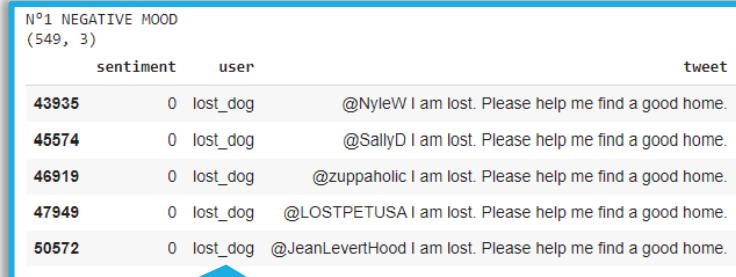


Déployer le meilleur modèle



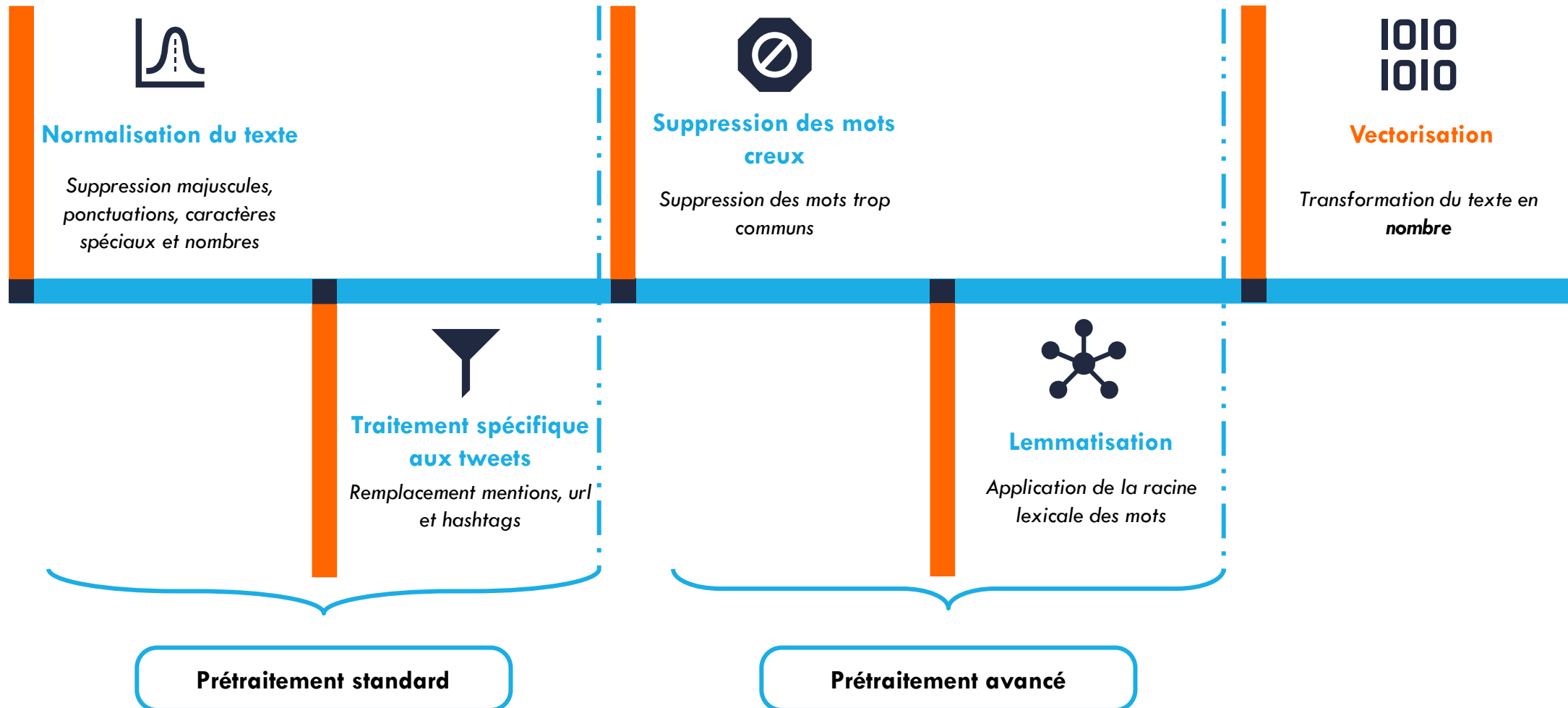
PRÉ-TRAITEMENT DES DONNÉES

Classes cible équilibrées



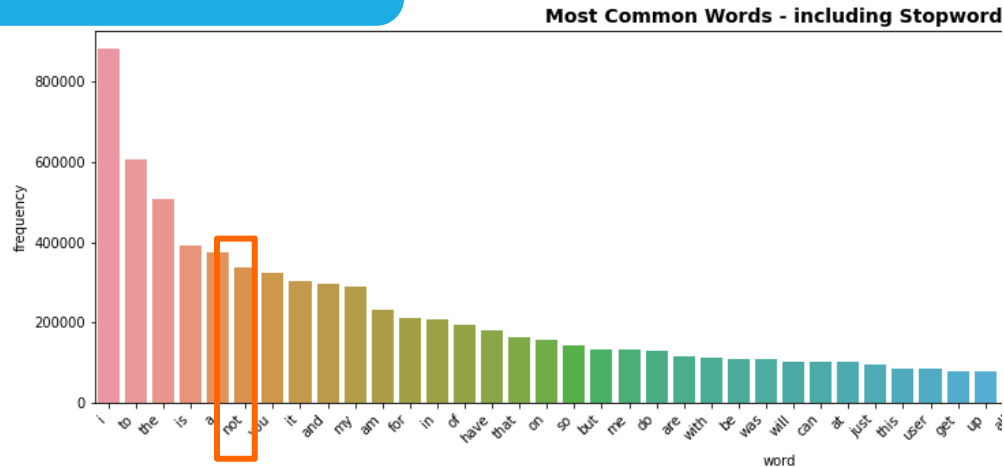
APPLICATION DE 2 TYPES DE PRÉTRAITEMENT

Nettoyage préliminaire : Suppression des doublons sur le sous-ensemble [Utilisateur – Tweet] et des lignes vides



COMPARAISON DES PRÉTRAITEMENTS

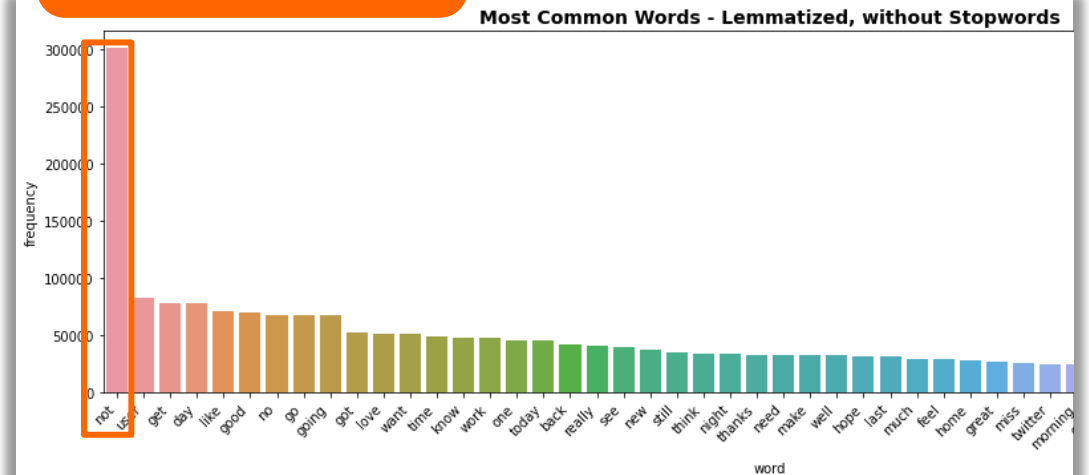
Prétraitement standard



'user nope they did not have it'

'user no it is not behaving at all i
am mad why am i here because i
can not see you all over there'

Prétraitement avancé



user nope not

'user no not behaving mad not see'

Utilisation des données avec **Prétraitement standard** : garder la séquence des mots, qui est importante dans l'analyse de sentiment



PRÉSENTATION DES APPROCHES

3 APPROCHES DIFFÉRENTES



Modèle « Prêt-à-l'emploi »



Bibliothèque de modules pour
construire des modèles



TensorFlow

Modèles « propriétaire » de réseaux
de neurones profonds

Choix de taille d'échantillon sur les tweets pré-traités pour un prototypage rapide

Entraînement et prédiction sur environ 10min par modèle

1.000 tweets

30.000 tweets

400.000 tweets

*Vectorisation si applicable

*Entraînement si applicable

Hachage des
caractéristiques

Extraction des
caractéristiques
n-grams

Word Embeddings
(plongement de mots)

Régression logistique

Réseaux de neurones récurrents

Prédiction

Evaluation des 3 approches + Déploiement du meilleur modèle

MODÈLE CLÉ-EN-MAIN : API SENTIMENT ANALYSIS DE MICROSOFT AZURE



Azure Cognitive Services

Prérequis:

Ouverture d'un compte Azure (sur le Portail)

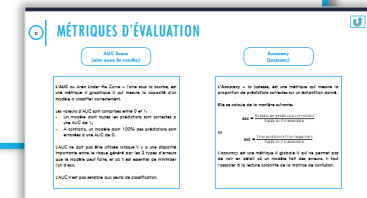
Souscription au service **Text Analytics** de la collection Azure Cognitive Services

Récupération du point de terminaison et de la clé de souscription pour utiliser le service

Utilisation du Kit de Développement Logiciel (SDK) Azure Machine Learning pour Python pour appeler le service

Pour valider notre hypothèse sur le prétraitement des données, nous avons testé l'API (Interface de programmation d'application avec les données :

- **Brutes**, sans aucun traitement;
- **Prétraitées**, standard;
- **Prétraitées**, avancé.



Model	Predict_time	AUC_Score	Accuracy
Original tweets	212.5	71.081%	70.860%
Cleaned tweets	197.2	76.187%	76.148%
Lemmatized tweets	198.2	74.595%	74.633%

MODÈLE BOÎTE-À-OUTILS : AZURE MACHINE LEARNING STUDIO (AMLS)



AZURE Machine Learning

Prérequis:

Ouverture d'un compte Azure (sur le Portail)
Souscription d'un espace de travail **Azure Machine Learning**

Clonage d'un pipeline existant

Microsoft Azure Machine Learning

Home > Designer

Designer

New pipeline

Easy-to-use prebuilt modules

Image Classification using DenseNet

Binary Classification using Vowpal Wabbit Model - Adult Income

Regression - Automobile Price Prediction (Basic)

Binary Classification with custom Python script - Credit Default

Binary Classification - Customer Relationship Prediction

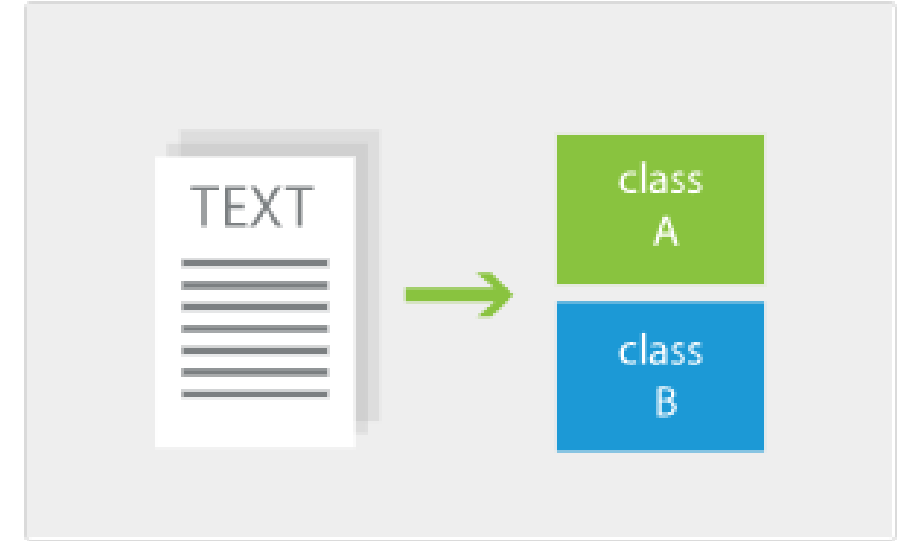
Use custom R script - Financial Delay Prediction

Text Classification - Wikipedia SP 500 Dataset

Cross Validation for Binary Classification - Adult Income

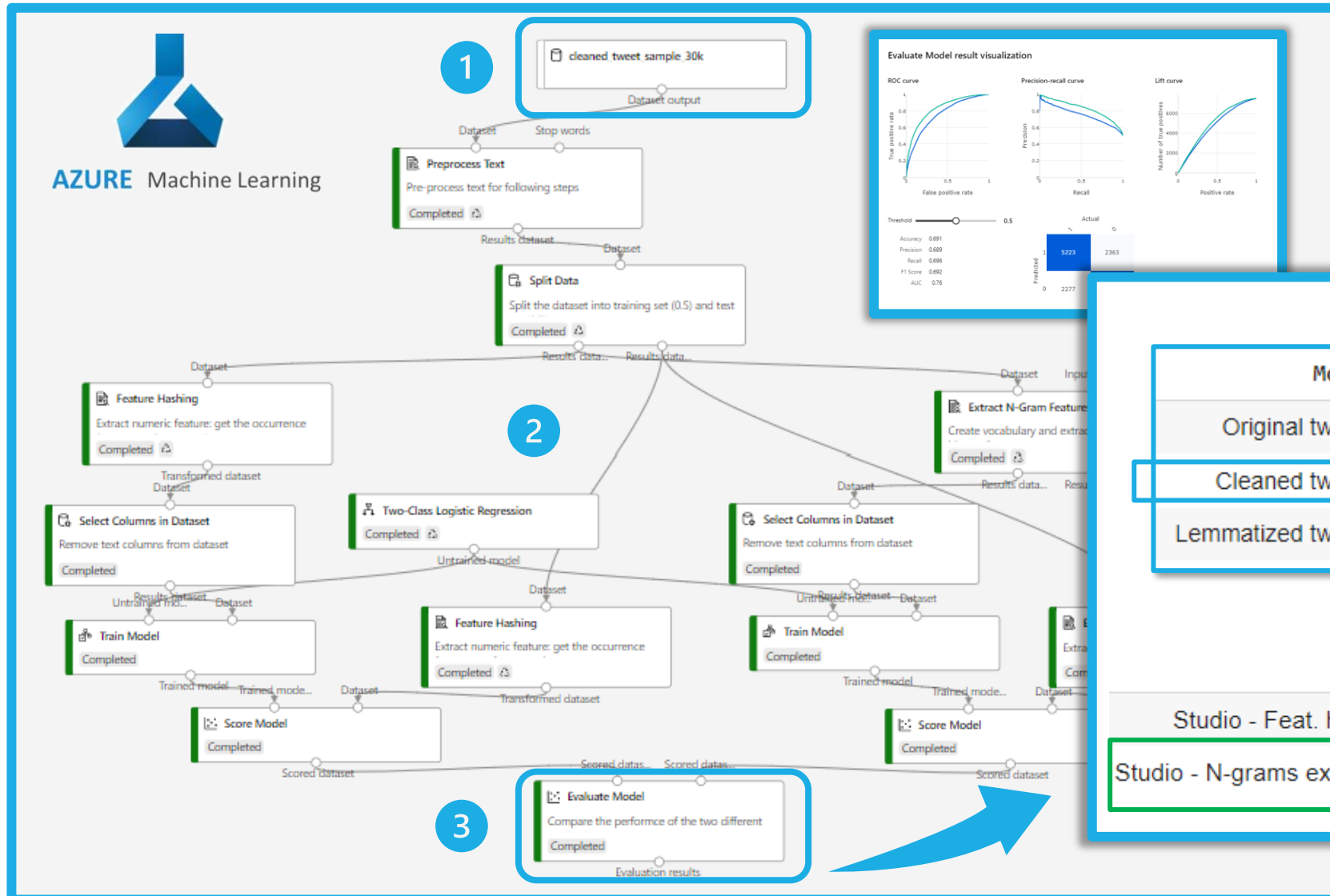
Tune Parameters for Binary Classification - Adult Income

Multiclass Classification - Letter Recognition



Text Classification - Wikipedia
SP 500 Dataset ?

AMLS : CONCEPTEUR ET PIPELINE DE MODULES



- 1 Remplacement du jeu de données
- 2 Correction des paramètres des modules
- 3 Exécution et évaluation des modèles

Approche 1

Model	Predict_time	AUC_Score	Accuracy
Original tweets	212.5	71.081%	70.860%
Cleaned tweets	197.2	76.187%	76.148%
Lemmatized tweets	198.2	74.595%	74.633%

Approche 2

Model	Predict_time	AUC_Score	Accuracy
Studio - Feat. hashing	81.5	76.900%	70.000%
Studio - N-grams extraction	103.4	84.200%	76.300%

MODÈLE AVANCÉ: RÉSEAUX DE NEURONES PROFONDS



Optimisation du modèle de base (baseline)

```
1 # Add sequential model
2 base_model = Sequential()
3 base_model.add(Embedding(input_dim=VOCAB_SIZE,
4                           output_dim=EMBEDDING_DIM,
5                           input_length=MAX_LEN))
6 base_model.add(SimpleRNN(128, dropout=0.3))
7 base_model.add(Dense(64, activation='relu'))
8 base_model.add(Dense(1, activation='sigmoid'))
9 base_model.summary()
```



Ajout de Dropout (0,3) et Regularizers (L2)

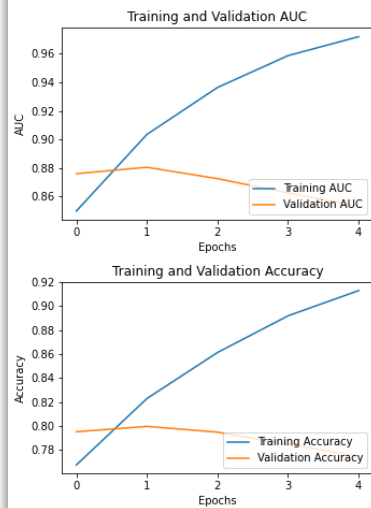
```
1 # Add sequential model
2 do3l2_model = Sequential()
3 do3l2_model.add(Embedding(input_dim=VOCAB_SIZE,
4                           output_dim=EMBEDDING_DIM,
5                           input_length=MAX_LEN))
6 do3l2_model.add(SpatialDropout1D(0.3))
7 do3l2_model.add(SimpleRNN(128, dropout=0.3))
8 do3l2_model.add(Dropout(0.3))
9 do3l2_model.add(Dense(64, activation='relu',
10                       kernel_regularizer=regularizers.l2(0.001)))
11 do3l2_model.add(Dropout(0.3))
12 do3l2_model.add(Dense(1, activation='sigmoid'))
13 do3l2_model.summary()
```



Optimisation du Dropout à 0,5

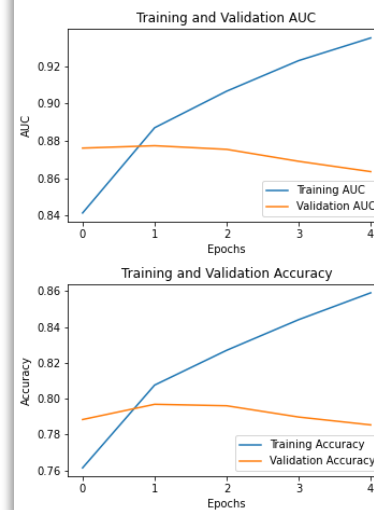
```
1 # Add sequential model
2 do5l2_model = Sequential()
3 do5l2_model.add(Embedding(input_dim=VOCAB_SIZE,
4                           output_dim=EMBEDDING_DIM,
5                           input_length=MAX_LEN))
6 do5l2_model.add(SpatialDropout1D(0.5))
7 do5l2_model.add(SimpleRNN(128, dropout=0.5))
8 do5l2_model.add(Dropout(0.5))
9 do5l2_model.add(Dense(64, activation='relu',
10                       kernel_regularizer=regularizers.l2(0.001)))
11 do5l2_model.add(Dropout(0.5))
12 do5l2_model.add(Dense(1, activation='sigmoid'))
13 do5l2_model.summary()
```

Base model RNN Model Learning Curves



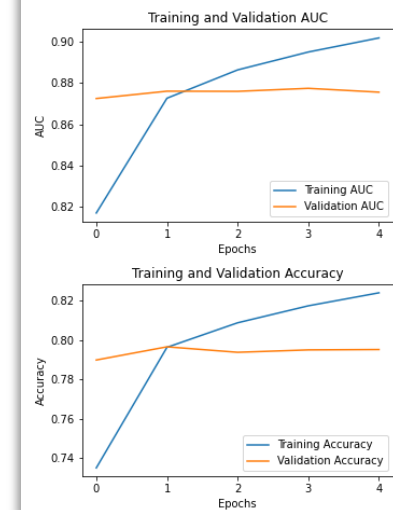
	Model	AUC_Score	Accuracy	Loss
0	SimpleRNN	85.245085	77.525002	56.828576

do3L2 baseline learning curves



	Model	AUC_Score	Accuracy	Loss
0	do3L2 SimpleRNN	86.35236	78.49375	49.211758

do5L2 baseline learning curves



	Model	AUC_Score	Accuracy	Loss
0	do5L2 SimpleRNN	87.459093	79.468751	44.949293

MODÈLE AVANCÉ: DEEP LEARNING & PLONGEMENT DE MOTS



TensorFlow

Paramètres communs

Configure all modeling variables

VOCAB_SIZE: 50000

MAX_LEN: 36

EMBEDDING_DIM: 256

DROP_OUT: 0.5

OPTIM_LR: 0.001

REGUL_LR: 0.001

NUM_EPOCHS: 5

BATCH_SIZE: 250

Modèle de base qu'on veut encore améliorer

Model: "base_model"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 36, 256)	12800000
spatial_dropout1d (SpatialDr	(None, 36, 256)	0
simple_rnn (SimpleRNN)	(None, 128)	49280
dropout (Dropout)	(None, 128)	0
dense (Dense)	(None, 64)	8256
dropout_1 (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 1)	65
Total params: 12,857,601		
Trainable params: 12,857,601		
Non-trainable params: 0		

Meilleure mémoire du passé proche

lstm (LSTM)	(None, 128)	197120
-------------	-------------	--------

Lecture gauche/droite et droite/gauche

bidirectional (Bidirectional	(None, 256)	394240
------------------------------	-------------	--------

Utilisation du modèle LSTM
avec les Embeddings pré-entraînés

Word2Vec

embedding_3 (Embedding)	(None, 36, 300)	28044000
-------------------------	-----------------	----------

GloVe

embedding_4 (Embedding)	(None, 36, 200)	18696000
-------------------------	-----------------	----------

USE

keras_layer (KerasLayer)	(None, 512)	256797824
--------------------------	-------------	-----------

Transforme une phrase entière en vecteurs, n'utilise pas de LSTM

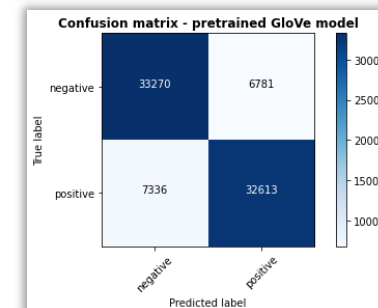
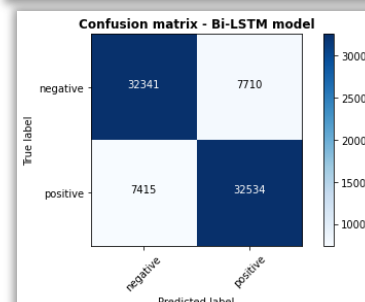
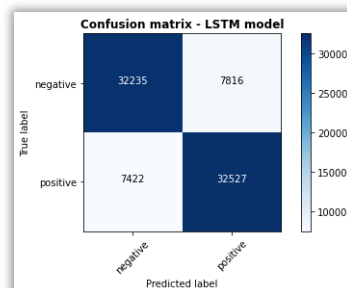
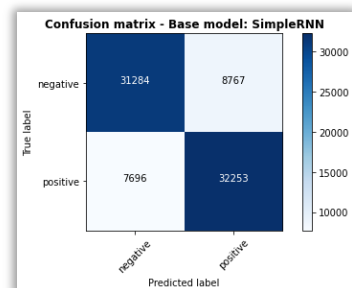
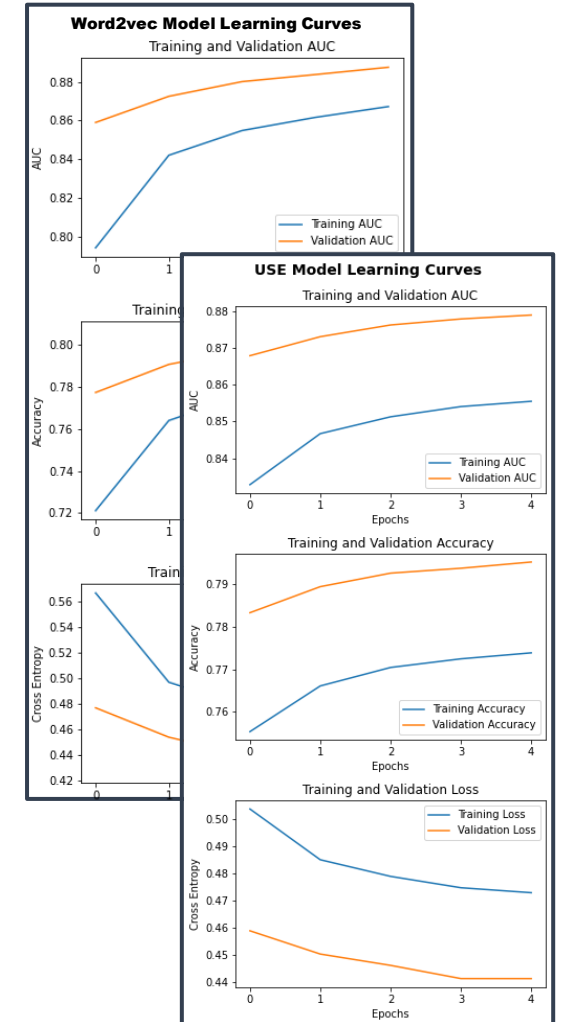
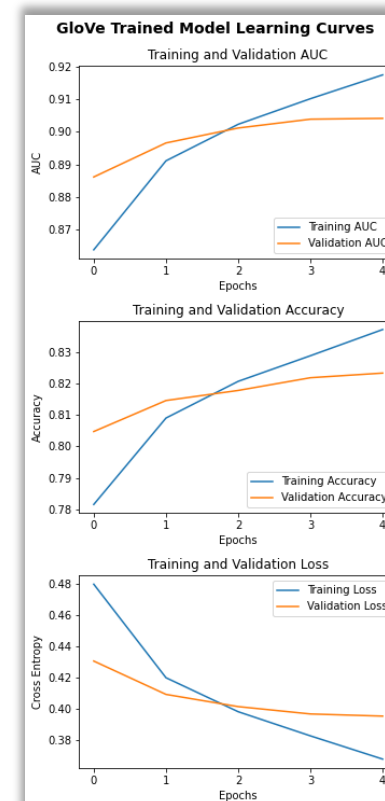
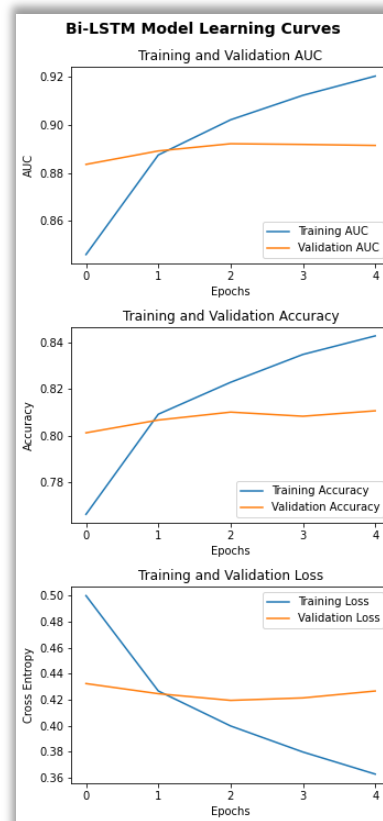
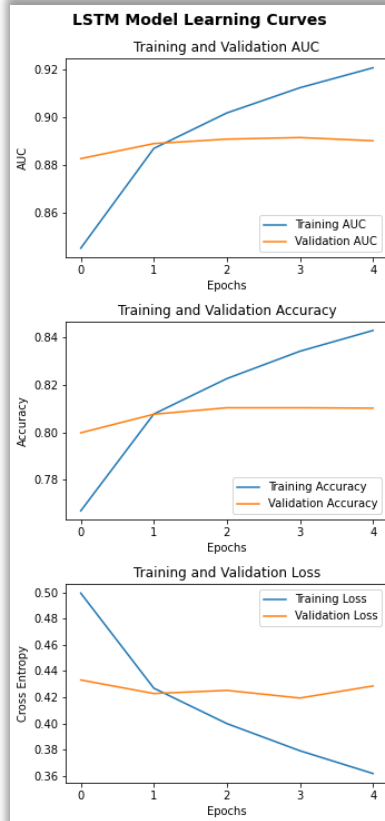
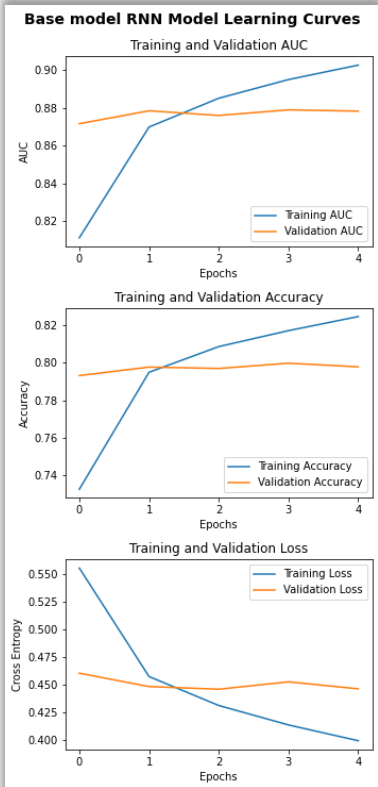
Différents réseaux de neurones : RNN → LSTM → LSTM bidirectionnel

Différents embeddings : from scratch → pré-entraînés (Word2Vec, GloVe, USE...)

PERFORMANCE DES MODÈLES AVANCÉS



Word2Vec et USE en sous-apprentissage (underfit)





CHOIX MODÈLE ET DÉPLOIEMENT

SYNTHÈSE DES PERFORMANCES DES APPROCHES

	Vectorisation	Modèle	Machine	Dimension de vecteurs	Echantillon	Paramètres entraînaibles	Temps d'exécution	AUC Score	Accuracy
	N.A.	API Sentiment	CPU	N.A.	1.000	N.A.	3,28 min	76,187 %	76,148 %
	Hachage de caractéristiques	Régression logistique	vCPU	N.A.	30.000	N.A.	~8 min	76,900 %	70,000 %
	Extraction de caract. n-grams	Régression logistique	vCPU	N.A.	30.000	N.A.	~7,30 min	84,200 %	76,300 %
	Embeddings propriétaire (Keras)	RNN Simple	GPU	256	400.000	12.857.601	8,10 min	87,653 %	79,421 %
	Embeddings propriétaire (Keras)	LSTM	GPU	256	400.000	13.005.441	7,41 min	88,955 %	80,953 %
	Embeddings propriétaire (Keras)	LSTM bidirectionnel	GPU	256	400.000	13.210.753	7,23 min	89,091 %	81,094 %
	Embeddings pré-entraînés: Word2Vec	LSTM	GPU	300	400.000	227.969	1,20 min	88,685 %	80,586 %
	Embeddings pré-entraînés: GloVe	LSTM	GPU	200	400.000	176.769	1,04 min	90,395 %	82,354 %
	Embeddings pré-entraînés: USE	N.A.	GPU	512	400.000	131,585	4,16 min	87,912 %	79,594 %

EXEMPLES D'ERREURS DU MEILLEUR MODÈLE : TWEETS LABELISÉS « NÉGATIF » PRÉDITS « POSITIF »



Housemate has told me to shut up with the hysterical laughing.

Holiday Was FANTASTIC

Sugarsnap peas, carrots and califlower. what a lunch!

kind of very pissed off that my parentals are going to NYC in Septmember without me. but now I just figured out what they can buy me!

Rona Howarda? S tim to jde z kopce. Kdyz vezmu Beautiful Mind ... Sifra ... Andele. Od deviti k peti



housemate has told me to shut up with the hysterical laughing

holiday was fantastic

sugarsnap peas carrots and califlower what a lunch

kind of very pissed off that my parentals are going to nyc in septmember without me but now i just figured out what they can buy me

user rona howarda s tim to jde z kopce kdyz vezmu beautiful mind sifra andele od deviti k peti

INTERPRETATION

Erreur réelle du modèle

Erreur de label à l'origine? Ce tweet semble très positif

Détection de sarcasme compliquée

Mix d'émotions: joie, peine...

Tweets qui sont dans une langue autre que l'Anglais (ici, le tchèque)

DÉPLOIEMENT SUR AZURE MACHINE LEARNING

Flux de travail

- Installation de Azure SDK pour Python
- Inscription du modèle sur Azure
- Préparation d'un script pour initier le déploiement et prédire avec le modèle
- Définition des configurations: environnement, dépendances, inférence
- Exécution du déploiement
- Utilisation du service web déployé

The screenshot displays the Microsoft Azure Machine Learning portal interface. The main view shows the details of a deployed service named 'text-sentiment-service'. The 'Attributes' section lists various properties:

- Service ID: text-sentiment-service
- Description: --
- Deployment state: Healthy
- Compute type: Container instance
- Created by: Voahangy Joan Aleonard
- Model ID: tweet_sentiment_glove_lstm:1
- Created on: 5/6/2021 10:40:40 PM
- Last updated on: 5/6/2021 10:40:40 PM
- Image ID: --
- REST endpoint: <http://7d309052-5bd2-48a1-ad65-48689281283e.francecentral2.azureml.ms/text-sentiment-service>
- Key-based authentication enabled: false
- Swagger URI: --
- CPU: 0.5
- Memory: 3 GB
- Application Insights enabled: false

The 'Test' tab is selected, showing the 'Input data to test real-time endpoint' section. A text input field contains the JSON payload: `{"text": "user that is a bummer url hashtag"}`. A 'Test' button is visible next to the input field.

The 'Test result' section shows the output of the inference, with a 'parsed' button selected. The result is a JSON object:

```
{
  "label": "NEGATIVE",
  "score": 0.035684049129486084,
  "elapsed_time": 0.07053208351135254
}
```



SYNTHÈSE

SYNTHÈSE DES INVESTIGATIONS

Flux d'un projet IA (vision 360°)



Les promesses du Cloud

- Facilité de création de flux de travail en Machine Learning
- Mise à l'échelle (scale up / down) des services en fonction des besoins
- Disponibilité permanente des services
- Affranchissement des coûts élevés de machines
- Paiement à l'usage
- Accès aux dernières mises à jour, innovations et technologies

Focus sur les objectifs

- **Vision d'ensemble** du projet de Machine Learning : place du projet dans la stratégie globale, usage, délai de mise en marché, ...
- Prototypage en local pour contrôler la **faisabilité** du projet
- Analyse de toutes les **options** disponibles
- Vérification de la capacité (interne) à **implémenter et à manager** les services dans le cloud

Approche 1

Outil clé en main, idéal pour commencer (petit)



Approche 2

Utilisation de services managés afin de profiter de l'existence d'algorithmes déjà éprouvés



Approche 3

Utilisation de modèles propriétaires avancés à déployer dans le Cloud

[illegible]

*Le coût d'un GPU est moindre par rapport à un vrai coût matériel

CONCLUSION

Un **choix d'approche** dépendant principalement des besoins exprimés par les équipes en interne ET des ressources disponibles (humain, financier, temps).

Une **performance du modèle** dépendant principalement :

- La qualité initiale des données ;
- Du prétraitement effectué.

The background of the slide is a complex network diagram. It consists of numerous small grey dots connected by thin grey lines, forming a web-like structure. Several larger circles are also present: a large white circle with a dark blue center at the top, a large blue circle at the bottom left, and a large grey circle at the bottom center. There are also smaller blue and dark blue circles scattered throughout the network.

QUESTIONS / RÉPONSES





ANNEXES

VECTORISATION DANS AZURE MACHINE LEARNING STUDIO (AMLS)



La **vectorisation** est le fait de transformer du texte en nombre pour qu'il puisse être ingérer par un modèle de Machine Learning.

Hachage de caractéristiques

Le **Hachage de caractéristiques** transforme un flux de texte en un ensemble de caractéristiques numériques, similaire à du cryptage. Un espace de hachage de taille fixe est définie (par exemple, 1.024 – et autant de colonnes) et le texte est converti en une séquence d'index dans cet espace.

Feature Hashing result visualization

Rows 24 000 Columns 1 027

tweet	sentiment	Preprocessed clean_tweet	Preprocessed clean_tweet_HashingFeature .0	Preprocessed clean_tweet_HashingFeature .1	Preprocessed clean_tweet_HashingFeature .2	Preprocessed clean_tweet_HashingFeature .3
ow we are o almost	1	me know we are open to almost anything	0	0	0	0
there is a devon love	1	seems there is a lot of devon love on twitter !!!	0	0	0	0
ter guess ot in id	1	guess i am not in splendid isolation after all	0	0	0	0
on after all re guy	1	user the guy next to you has what can only be described as a worrying expression	0	0	0	0.174078
y you has an only be ned as a	1	ng	0	0	0	0
sion et on it way	0	I will get on it right away	0	0	0	0
ay is the at people	0	so today is the day that people are lining up for the iphone gs !!!	0	0	0.140028	0
ing up for the phone gs i	0	can not get the discount until	0	0	0	0
i get the nt until	0		0	0	0	0

Extraction de caractéristiques n-grams

L'**extraction de caractéristiques de n-grams du texte** extrait d'un flux de texte un vocabulaire, affecte un valeur de présence (1) aux n-grams extraits lorsqu'ils existent dans le document, et attribue un score de fréquence (TF-IDF) à ces mêmes n-grams.

Extract N-Gram Features from Text result visualization

Results dataset Result vocabulary

Rows 24 000 Columns 8 391

clean_tweet	sentiment	Preprocessed clean_tweet	Preprocessed clean_tweet. [able_watch]	Preprocessed clean_tweet. [about_going]
bad headache will not go away	0	bad headache will not go away	0	0
twilight getting so many awards in the mtv movie awards god i	0	twilight getting so many awards in the mtv movie awards !!! god i	0	0

Extract N-Gram Features from Text result visualization

Results dataset Result vocabulary

Rows 8 389 Columns 4

Id	NGram	DF	IDF
0	able_watch	5	3.681241
1	about_going	5	3.681241
2	about_too	5	3.681241
3	about_twitter	5	3.681241
4	about_user	5	3.681241
5	abuse	5	3.681241
6	accent	5	3.681241
7	addad_van	5	3.681241

WORDS EMBEDDINGS OU PLONGEMENT DE MOTS



Le **plongement de mot** est une méthode de structuration (vectorisation) des données textuelles.

L'idée est de représenter chaque mot issu d'un corpus par un **vecteur (liste de nombres) de dimension fixe**, qui va en quelque sorte « capturer son sens ».

L'hypothèse est la suivante : **les mots qui apparaissent souvent avec le même « contexte » (mots avant/après) sont probablement proches d'un point de vue sémantique, et peu distant dans un espace vectoriel.**

Schématiquement, si on considère les phrases suivantes :

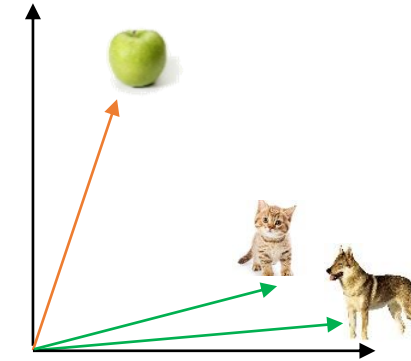
'Le chien est un animal'

'Le chat est un animal'

'La pomme est un fruit'

, on sait que *chien* et *chat* sont sémantiquement plus proches que *chien* et *pomme*, ou *chat* et *pomme*.

En présentant les vecteurs de ces trois mots sur 2 dimensions, on devrait voir la figure ci-après.



Il existe 2 principales méthodes d'entraînement des **word embeddings** :

- **CBOW** (Continuous Bag of Words), qui prédit un mot à l'aide du contexte (mots avant/après) ;
Ex: Dans : 'Le ____ miaule', on devine le mot *chat*
- **Skip-gram**, qui prédit – à contrario, le contexte (mots avant/après) à l'aide d'un mot ; Ex:
 '____ gros **chien** jappe'
 'Le ____ **chien** jappe'
 'Le gros **chien** ____'

WORDS EMBEDDINGS : À ENTRAÎNER OU PRÉ-ENTRAÎNÉS

Il existe 2 sortes de **Word Embeddings** :

- **Les embeddings « à entraîner »** : ils sont entraînés en temps réel sur les données d'entrée; on ajoute une couche d'Embedding à notre modèle, qui transforme les mots en vecteurs dans une matrice, celle-ci étant « apprise » au fur et à mesure que le modèle s'entraîne ;
- **Les embeddings « pré-entraînés »** : ils ont été pré-entraînés sur des corpus de texte très importants (Wikipedia, Google News, ...) et ils permettent un temps d'entraînement réduit et souvent, de meilleures performances globales.

Embeddings	Editeur	Année	Niveau	Téléchargement	Type	Architecture	Dimensions
Word2Vec	Mikolov et al	2013	Mots	~1.5GB	À entraîner	CBOW et SG	300
GloVe* (Global Vector for word representation)	Pennington et al	2014	Mots	~1.5GB	Pré-entraîné	Ratio de cooccurrence	25, 50, 100, 200
USE (Universal Sentence Encoder)	Google	2018	Phrases	~1GB	Pré-entraîné	DAN (Deep averaging network)	512

*Version Twitter

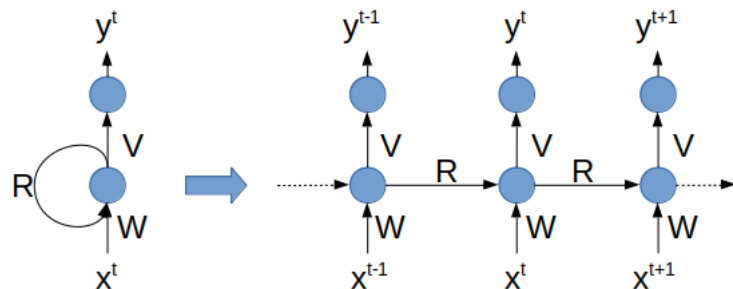
RÉSEAUX DE NEURONES RÉCURRENTS (RNN)



Les **réseaux de neurones récurrents (ou RNNs)** sont un type de réseaux de neurones artificiels conçus de manière à **reconnaître les caractéristiques séquentielles** telles que les séries temporelles ou le langage naturel (textuel ou vocal).

En cela, ils sont plus proches du fonctionnement réel du système nerveux humain, au sens où ils **conservent en mémoire** un certain nombre d'états passés à un instant t .

Schématiquement, on dit que les RNNs utilisent des **boucles de rétroaction** pour traiter une séquence de données, afin de permettre aux informations de « persister » dans le temps.



RNN : vision récurrence vs. vision déplié (temporalité)

Source: [Open Classrooms](#)

Les RNNs utilisent la méthode d'apprentissage par **rétropropagation du gradient à travers le temps** pour apprendre les paramètres du réseaux. Néanmoins, ceux-ci tendent à se dégrader à mesure qu'ils grossissent et se complexifient.

Il existe différents RNNs selon l'étendue de leur « mémoire à court-terme »:

- Les **RNNs** « simples », qui oublient au bout de quelques itérations le passé dit proche ;
- Les unités **LSTMs** (Long Short-Term Memory), qui ont à la fois une mémoire long et court-terme, donc avec la capacité de distinguer les données importantes à mémoriser pour les réinjecter dans le réseau ; elles sont 2 fois plus lourdes que les couches récurrentes simples, compte tenu des connexions nécessaires.

Des variantes sont également disponibles pour les LSTMs:

- Les GRU (Gated Recurrent Units), qui réduisent le nombre de paramètres des modèles ;
- Les LSTM bidirectionnels (BLSTM), qui observent le passé mais également le futur.

MÉTRIQUES D'ÉVALUATION



AUC Score (aire sous la courbe)

L'AUC ou *Area Under the Curve* – l'aire sous la courbe, est une métrique « graphique » qui mesure la capacité d'un modèle à classifier correctement.

Les valeurs d'AUC sont comprises entre 0 et 1 :

- Un modèle dont toutes les prédictions sont correctes a une AUC de 1 ;
- A contrario, un modèle dont 100% des prédictions sont erronées a une AUC de 0.

L'AUC ne doit pas être utilisée lorsque il y a une disparité importante entre le risque généré par les 2 types d'erreurs que le modèle peut faire, et où il est essentiel de minimiser l'un d'eux.

L'AUC n'est pas sensible aux seuils de classification.

Accuracy (justesse)

L'Accuracy – la justesse, est une métrique qui mesure la proportion de prédictions correctes sur un échantillon donné.

Elle se calcule de la manière suivante:

$$\text{acc} = \frac{\text{Nombre de prédictions correctes}}{\text{Taille de l'échantillon}}$$

ou

$$\text{acc} = \frac{\text{True positives} + \text{True negatives}}{\text{Taille de l'échantillon}}$$

L'accuracy est une métrique « globale » qui ne permet pas de voir en détail où un modèle fait des erreurs. Il faut l'associer à la lecture conjointe de la matrice de confusion.

TARIFICATION AZURE PAR SERVICE



Cognitive Services <clé en main> Text Analytics > Sentiment

Gratuit

5000 transactions gratuites par mois

Standard (dégressif)

Tous les 1.000 enregistrements de texte

0 à 500k enregistrements	0,844€
500k à 2,5 millions	0,633€
2,5 à 10,0 millions	0,253€
> 10 millions	0,211€

Instance Clusters de calcul Azure ML Studio (apprentissage)

Virtual Machine CPU

Options recommandées:
Machine virtuelle dédiée,
de type CPU (vCPU),
avec un quota de 20 cœurs
(coût par heure)

Usage général : D2 – 2 cœurs – 7GB ram – 14Go stockage	0,148€
Usage général : D3 – 4 cœurs – 14GB ram – 28Go stockage	0,296€
Mémoire optimisée : DS12 – 4 cœurs – 28GB ram – 28Go stckg	0,396€
Calcul optimisé :F4S – 4 cœurs – 8GB ram – 32Go stockage	0,396€

Clusters d'inférence Azure ML Studio (déploiement)

Création de conteneur (€/UC)

0,00009€/seconde de tâches en cours
d'exécution pour Basic, Standard et Premium

0,00009€/seconde d'allocation d'instance
(nombre de processeurs alloués) pour les Pool
d'agents dédiés: S1-2c-3Go ram, S2-4c-8Go
ram ou S3-8c-16Go ram

Registre de conteneurs

Tarif quotidien + Tarif par région répliquée

Basic 10Go stockage	0,141€
Standard 100Go stockage	0,563€
Premium 500Go stockage Géoréplication (/région répliquée)	1,406€ 1,406€

Stockage additionnel: 0,003€/jour

RÉFÉRENCES

- ❑ Analyse de texte (text mining) avec [Azure Cognitive Services - Text Analytics](#)
- ❑ Solution Cloud de Machine Learning avec [Azure Machine Learning](#) et son [concepteur \(Drag and Drop\)](#)
- ❑ Deep Learning : [Tensorflow - Réseaux de neurones récurrents](#), [LSTM - Towards Data Science](#) , [word embeddings](#)
- ❑ Pretrained word embeddings : [word2vec](#), [GloVe \(Global Vectors for Word Representation\)](#) , [USE \(Universal Sentence Encoder\) - v4](#)
- ❑ Déploiement : [How and where to deploy Machine Learning models to Azure](#)
- ❑ Autre : [neptune.ai - structure and manage nlp projects](#)
- ❑ Pour aller plus loin : [Classification de texte avec BERT](#), [Traduction de texte avec des modèles seq2seq \(NMT: Neural machine translation with attention\)](#) , [Traduction de texte avec des modèles Transformer](#)



Ce document a été produit dans le cadre de la soutenance du projet n°7 du parcours Ingénieur IA d'OpenClassrooms :
« Détectez les Bad Buzz grâce au Deep Learning »

Mentor : Thierno DIOP
Evalueur : Kezhan SHI

