

Comparative Study To Solve Text Categorization

CES Data Scientist 2017-2018 (Telecom ParisTech)

Jacques Doan-Huu

Abstract

In the last decade, **Deep Learning (DL)** has demonstrated outstanding performance in the field of computer vision beating indisputably traditional methods, thanks to the GPU performance leapfrog and the huge amount of labeled datasets. A bit more recently, DL also went into the **Natural Language Processing (NLP)** field battle to solve common NLP problems like text classification and translation, with very promising perspectives and results: in particular, word embedding and **Recurrent/Convolutional Neural Network (RNN/CNN)** architectures provide efficient technical responses to NLP challenge.

POSOS French startup has submitted a data challenge for which I took the opportunity to verify humbly whether DL is a suitable solution compared to traditional methods, for a beginner like me having very few practices on NLP/DL area and low-end hardware (DL has the bad reputation to be numerically intensive...).

Contents

Abstract	1
Statement of the Problem	3
Project Motivation	3
ML Workbench Environment.....	3
Data Exploration	4
Target Distribution	4
Topic Extraction	4
Feature Space Distribution	5
Text Anatomy Analysis.....	7
General NLP Architecture	8
Text Preprocessing.....	9
Tokenization	9
Spelling correction	9
Lexical and Grammar Tagging.....	11
Text Cleansing and Normalization	12
Classical Technique	13
Abstract.....	13
Feature Enrichment	13
Feature Representation	15
Classification Modeling.....	16
Architecture	16

Hyperparameter Search.....	16
Result Analysis	17
Deep Learning Technique	19
Abstract.....	19
Feature Enrichment	19
Feature Representation	19
Sequential representation	19
Word embedding.....	19
DNN Architecture.....	21
CNN Architecture	21
RNN Architecture	22
Hyperparameter search and Architecture Sizing.....	24
DNN.....	24
Sequence Padding.....	24
Parallel CNN	24
Sequential CNN	24
RNN / LSTM.....	24
DL Implementation and Execution	25
Result Analysis	26
Comparative Study	27
Model Accuracy	28
Model Interpretability	28
Tooling	29
Sustainability.....	29
Improvement Tracks.....	29
Spelling Correction.....	30
Named Entity Recognition	30
Count/Distance based Statistics	30
Word Embedding	31
OOV Handling	31
Early Stopping With Grid Search CV.....	31
Neural Network Tuning.....	31
Other Modeling Candidates.....	32
Too Few Samples and Too Many Target Effects	32
Conclusion	34
Appendix.....	35
Github project.....	35
References	35
Resources (data/pretrained model)	35

NN Papers / Blogs	35
Tools.....	36

Statement of the Problem

The data challenge is plainly described at ENS school web site ([link](#)) and it consists in categorizing into **51** intents, drug related questions written in natural language (French to be precise). **POSOS** claimed to get good performance with 86% accuracy by choosing DL: they don't supply any details on the DL architecture nor any engineering clues except the recommendation to extract some key information procured by the French drug administration (**ANSM**).

The target categories (question intent) have been intentionally anonymized into indices from 0 to 50: hiding their respective semantic is probably aimed to avoid the usage of topic-specific (and so biased) procedures. Training dataset contains only ~8000 questions: it's pretty short to produce a good learning outcome. Besides, the text suffers from many anomalies (misspelling, grammatical incorrectness, familiar acronym, ...) and employs specific medical vocabulary (drug name like "mirtazapine", ...): it hardens the challenge level of difficulty.

Project Motivation

The purpose of this study is to compare fairly the strength and weakness between DL and non-DL approaches from different perspectives:

- model accuracy
- model interpretability
- tooling
- sustainability

This is not about achieving high score at any price, but an attempt to explore comparatively the end to end methodology to tackle a text categorization problem throughout 2 distinct technologies.

"Traditional techniques" refer to any ML algorithms which don't rely on neural network theory (eg: Word2Vec is excluded): to quote some of them, hidden Markov model, gradient boosting, non-negative matrix factorization, logistic regression and principal component analysis are eligible.

Conversely, DL option should be uniquely based on neural network but it can as well benefit from "neutral" text preprocessing (feature enrichment with external source, stemming, stopWords, ...) for fairness sake.

ML Workbench Environment

All experiments have been written in Python in the popular Jupyter environment: the notebooks are available publicly as a Github project whose details are provided in the appendix section. I used many python packages to satisfy various requirements:

- data manipulation and visualization: pandas, numpy, seaborn and matplotlib
- text processing (stemming, stopWords, ...): NLTK, standard regex and spellChecker (built from github)
- ML algorithms (XGBoost, PCA, t-SNE, NMF): sklearn and XGBoost
- DL framework: Keras + Tensorflow

Most of packages have been installed as is, except for XGBoost I recompiled locally from its github source code to get the GPU accelerated version which is not shipped officially.

Besides above runtime packages, this project also takes advantage of public resource or pretrained models (NLTK corpus, FastText word embedding model, ...).

ML jobs had been initially executed with an old MacBook Pro whose chipset was damaged by the heating due to the overnight DL train, then with a many CPU core/low-end GPU PC workstation.

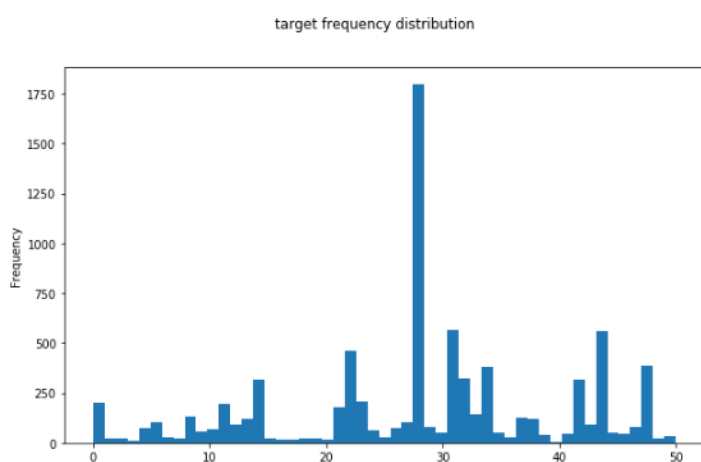
I finally rent a GPU boosted Amazon instance: even if the GPU/CPU resource is not utilized, the data storage is charged permanently making the overall cost very high (5\$/h) during 2 full days.

At the very end, the best money saving option was to buy a PC gamer machine with mid-range Nvidia GPU card to carry out the computing workload: I have gained an important speedup at training time.

Data Exploration

Target Distribution

The target distribution of the 51 labels is imbalance with a peak at intention=28 with 273 as standard deviation.



Most of labels are associated to pretty small number of samples: half of classes have less than 100 rows, it sounds that the classifier would underperform on such labels with low input features.

Topic Extraction

Let's try to guess intuitively the hidden meaning of the most frequent labels:

first mode: intention=28

It likely corresponds to drug adverse effects (contraindications)

0	bonjour, je m suis trompé de forum pour ma question alors je la repose ici. je pris pour la première fois hier du paroxétine et ce matin c'est une catastrophe. picotement dasn tous le corps annonciateur de sueur froide très très massive et de vomissement. j'en suis à deux crises depuis 5 heure du mat. la cela semble passer mes mes mains reste moites et chaude estce normal pour la première fois merci a tous	28
2	mon médecin m'a prescrit adenyli. au 2ème cachet des maux de tête terribles et au 3ème palpitations, sueurs froides, chaleur intense dans la tête, tremblements, fourmillements dans la lèvre supérieure, difficultés à respirer.. dès l'arrêt du médicament tous les symptômes ont disparu. cela est-il déjà arrivé à quelqu'un??	28
7	je suis sous mercilon. J'ai des nausées et des saignements ?	28
12	je suis sous antibiotique depuis bientôt une semaine et je me suis chopée je ne sais quoi à ma nénéte, ca gratte,c'est superficiel mais ca démange à un point, est ce lié à l'antibiotique?	28
14	épilepsie et havlane ?	28

second mode: intention=31

it's about drug>disease indication/efficiency

1	1	est ce que le motilium me soulagera contre les nausées?	31
4	4	mon medecin me soigne pour une rhino pharyngite et m'a prescrit du amoxicilline comme anti biotique. Est-ce vraiment pour cette indication?	31
10	10	laroxyl à doses faibles pour le stress ?	31
31	31	La lidocaïne aide-t-elle à maigrir ?	31
37	37	L'euphytose est utile pour l'anxiété ?	31

multi-topic class: intention=39

the commonality across text samples is the presence of multiple question marks (counting it may be a good option to predict multi-topic label)

on m'a prescrit microval, est ce que cette pilule est effective dès la première prise ? qu'est ce que cela fait si je commence la plaquette alors que je n'ai pas encore mes regles ?	39
mon medecin m'a prescrit du xanax 0,50 et du stilnox. cependant j'ai l'impression que le stilnox (generique) commence a ne plus me faire effet. faudrait t'il que je lui demande plutot de l'imovane et quels sont les effets secondaires?	39
PAS DE QUESTION SUR UN MEDICAMENT	39
samedi soir, j'ai oublié de prendre mon 17è comprimé et l'ai pris le lendemain mais pense l'avoir vomi. 4 jours avant l'oubli de ma pilule, j'ai eu un rapport non protégé : ai-je un risque de grossesse dû à ce rapport ? ma plaquette se finit dans 3 jours : au démarrage de la prochaine, puis-je recommencer à avoir des rapports non protégés?	39
ai arrêté mon traitement de dépanage il y a un mois et demi et malgré le sport et un régime sévère, je ne vois pas de changement sur ma balance.. est ce que certaines parmi vous ont eu ce pb ? combien de temps vous a-t-il fallu pour retrouver votre poids normal??	39

Instead of guessing manually the label meaning, I made use of **NMF** (Non-Negative Matrix Factorization) algorithm to extract the main topics: I fixed the number of topics to the number of labels in the naïve hope of finding an exact match.

Topic #0: depuis prend plus semain ça an normal quelqu mal cel tout pens problem peu ca arriv ventr tres cet comm do uleur bonjour saign merc pert
Topic #1: secondair effet infanrixquint tolexin zoloft influenzinum microval don trinordiol rotarix prescr lexiomil citalo pram gripp beaucoup abilify laroxyl zyprex dostinex foliqu lutéran lutenyl tamik minidril dompéridon
Topic #2: quel dosag dos posolog médic thérapeut altern indiqu maximal form leponex vitamin moment action class uti l recommand cas différent appartient rivotril arnic prix effet enfant
Topic #3: grossess pend possibl levothyrox prescrire dur champix aeriux dang flagyl test autoris ginkor essentiel début c ompatibl vogalen subutex lysopain azithromycin depakot semain efferalgan primperan danger
Topic #4: vaccin gripp hépatit dtp varicel ror gardasil exist rappel fievr inject hepatit tétanos polio col méningit contr a pres an réaction jaun utérus dt où fil
Topic #5: combien temp bout efficac effet faut dur attendr agir met agit apres fass mem xenical fait granul sertralin lum ali norlevo apre pend conserv durent sang

I observed that some topics actually correspond to some labels:

- topic #1 matches the label 28 (drug adverse effect)
- topic #3 matches the label on drug and pregnancy interaction
- topic #4 is about vaccine

Such identified topics will be used later on as an extra feature to improve the classification rate.

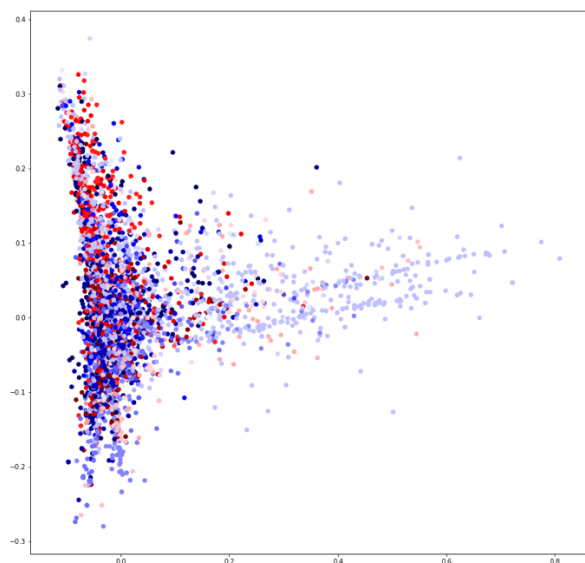
Feature Space Distribution

To have a rough estimate on the classification task difficulty, it's a common practice to visualize the feature space distribution.

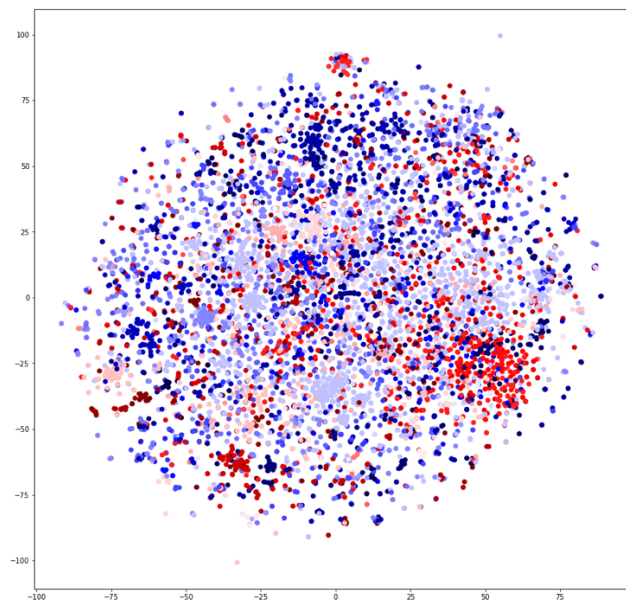
The document is basically transformed into of **BOW** (Bag Of Words) which is then vectorized with TF-IDF encoder: each document is consequently represented as a data point within the global vocabulary space. To make such data

points human readable, a dimension reduction of these features is necessary at the cost of some approximations: I used both linear/fast PCA dimension reduction and non-linear/slow t-SNE techniques. The data point color reveals the associated target label.

PCA-reduced feature space



t-SNE-reduced feature space



Both 2D distribution shapes are very dissimilar but they consistently tend to indicate that feature space cannot be partitioned per label easily just by considering the occurrence of words. To get a chance to achieve a better classification performance, it's obvious that the raw text has to be encoded more smartly into a suitable space where the semantic proximity between documents is prevailing.

Text Anatomy Analysis

Let's look at 2 samples having the same label (disease-drug adequacy):

"**épilepsie** et **havlane**?"

"mon medecin me soigne pour une **rhino pharingite** et m'a prescrit du **amoxicilline** comme anti biotique. Est-ce vraiment pour cette **indication**?"

Even if they share the same question topic, the writing styles are completely opposite: on one hand, a very concise expression putting the disease entity and the drug entity in an adversarial fashion, on the other hand, the second sample is more descriptive and spread over 2 sentences.

This example is a manifest of the stylish complexity of the human language to convey an idea and a topic!

The second sample has 2 sentences: the first one settles the question context and the current situation ("mon médecin me soigne..." whereas the second one raises the effective question ("Est-ce que....").

In multi-sentence documents, I generally observed this sequential structure: first the context setup, followed by the concrete question.

3 entities turn out to be salient here:

- the drug product name ("amoxicilline", "havlane")
- the disease ("épilepsie", "rhino-pharingite")
- the link entity between above ("et", "indication")

Identifying the first 2 entities is doable just based on lexical semantic domain: merely, build an exhaustive list of symbols related to drug product or disease. It's related to **NER** (Named Entity Recognition) task.

The last entity connecting the 2 other entities is much more difficult to locate: the entity semantic is contextual and depends on the presence of other entities and its relative position within the text grammatical structure. It means that the learning procedure should be a **sequence modeling**.

Other tokens (mon, médecin, me, soigne, ...) seem to be superfluous to extract the question intent: text processing to remove irrelevant words is highly recommended (**stopwords**, custom regular expression, ..)

Last but not the least, some documents are lexically and syntactically incorrect: words are misspelled especially when dealing with drug product names which are unfamiliar for most of non-professional persons. A **spelling correction** is required in the text preprocessing phase.

The misspelling count average per document is 1.1 and drug name is present in most of questions (0.89) but ingredient entity is barely mentioned (0.22).

	drug name count	active ingredient count	misspelling count
count	8028.000000	8028.000000	8028.000000
mean	0.896487	0.229571	1.110737
std	0.663804	0.506578	1.777682
min	0.000000	0.000000	0.000000
25%	1.000000	0.000000	0.000000
50%	1.000000	0.000000	1.000000
75%	1.000000	0.000000	1.000000
max	10.000000	6.000000	27.000000

Regarding statistics on text shows that on average, the document is concise (10 words, 69 characters)

```

: XTrain['text length'] = XTrain['question length']
XTrain['text length'].describe()

: count      8028.000000
  mean        69.238540
  std         89.399312
  min          4.000000
  25%         28.000000
  50%         42.000000
  75%         80.000000
  max        2744.000000
  Name: text length, dtype: float64

XTrain['word count'] = XTrain['question word count']
XTrain['word count'].describe()

: count      8028.000000
  mean        9.854758
  std        13.273684
  min         1.000000
  25%         4.000000
  50%         6.000000
  75%        11.000000
  max        408.000000
  Name: word count, dtype: float64

```

The vocabulary is significantly big (9378 words) and words are on average infrequent (8 occurrences). When inspecting the most frequent words, common-used but low informative terms rank first and only 2 medication related terms come out (pilule and vaccin).

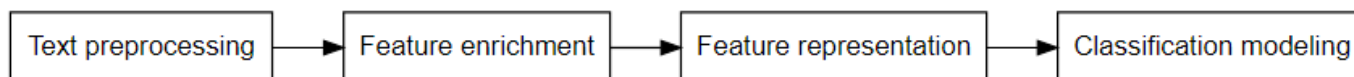
Another interesting point is the presence of 3 morphological variants of the root “prendre” in the top 20: a **stemming/lemmatization** or **word embedding** are welcome to collapse such semantically equivalent variants into a single representative.

word frequency	
count	9378.000000
mean	8.421305
std	42.488629
min	1.000000
25%	1.000000
50%	2.000000
75%	4.000000
max	1993.000000

word frequency	
a	1993
les	1389
si	827
depuis	780
peut	779
prendre	771
pilule	751
jours	593
plus	557
effets	554
faire	546
mois	529
sous	501
fait	470
savoir	468
mg	443
prends	431
temps	386
pris	380
vaccin	379

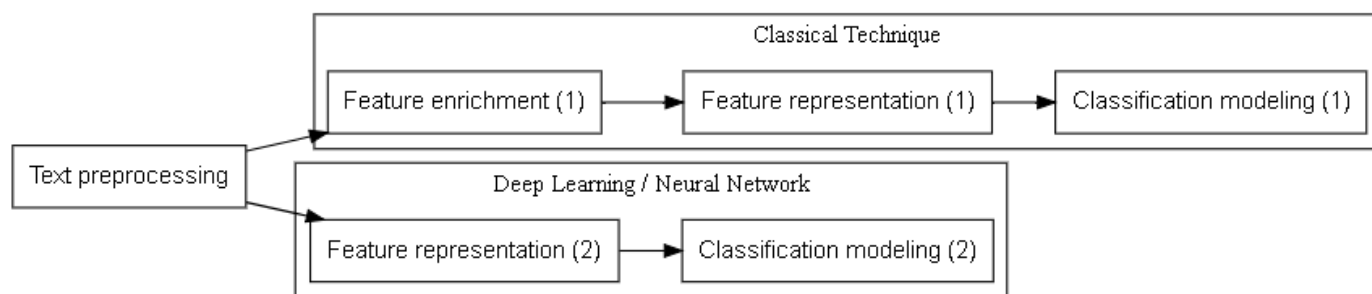
General NLP Architecture

Text classification is a standard but non-trivial NLP topic going through the following main steps:



This is the general NLP text classifier framework/guidance but for practical reasons, some processing steps are skipped or significantly simplified to fit the project timeframe but also because of the lack of French language support.

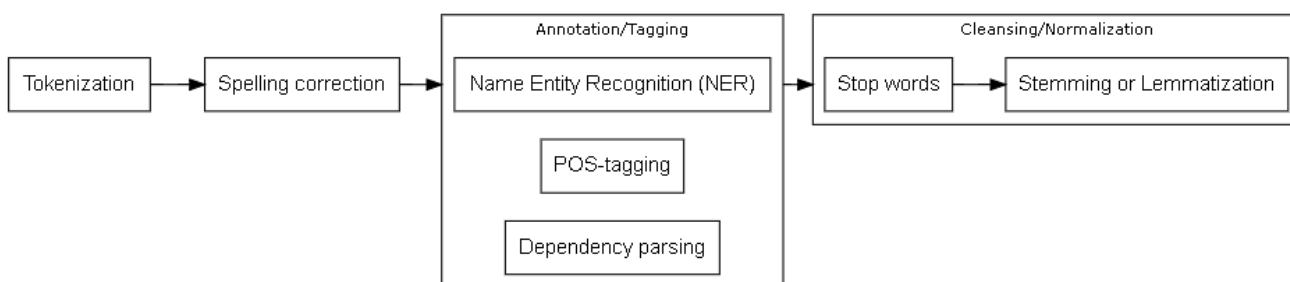
In fact, here’s the concrete pipeline I built per technical scenario:



Each processing unit will be described more precisely in the next sections.

The overall modeling procedure should capture the sequential nature of the text to exploit efficiently the contextual information: typically, the feature representation should preserve the word/symbol order and the classification process should be compliant with **sequence modeling**.

Text Preprocessing (*common trunk*)



It's all about operations on the raw text to make it more reliable/workable in order to extract relevant characteristics. It falls into 4 categories:

- tokenization breaking down the sentence into a sequence of atomic words/symbols
- correction on misspelled words
- lexical and grammar tagging which primarily decorates the text tokens with metadata
- text cleansing and normalization simplifying the sentence composition

Tokenization

This operation is a commonplace: I simply used the python string `split()` function. I tokenized the whole document by ignoring the punctuation like “.”, “:”, “;”, “!” and “?”.

Spelling correction

I assumed as misspelled all words which don't belong to any trustworthy vocabularies, also known as **OOV** (Out Of Vocabulary) word.

I retained 3 reference vocabularies in the following priority order:

- vocabulary from the word embedding model used downstream in the processing pipeline
 - indeed, it's very important to avoid random vectorization on OOV words

- Custom vocabulary to capture the specific drug domain, typically on drug product and active ingredient entities where misspelling is frequent. It has been built from the public RCP (Résumé des Caractéristiques du Produit) repository supplied by ANSM.
- Predefined general vocabulary from the github python project pysspellchecker
<https://github.com/barrust/pysspellchecker>

The curative algorithm finds from a set of vocabularies, the closest word candidate from **Levenhstein** distance standpoint: this distance measures the minimum number of character operations (change, remove, add) required between 2 words.

I defined an empirical threshold to accept the closest word as a fix on the misspelled word: the ratio between the number of atomic operations and the total number of characters should be under 25%.

I applied this algorithm with the last 2 vocabularies: the first vocabulary layer only filters out the recognized words, the unfixed words at the second layer are then passed to the third layer.

Here's the python output showing the fix on more than 400 drug product names with a reasonable error rate (~15%):

```
n [3]: drugNames = Counter(words(open('../data/staging_data/drug_names.1
missSpelledDrugMap = buildFixMap(unknownWords, drugNames, 0.25, None)

KO gényco => glyco with dist=0.3333
OK thyrosine => thyroxine with dist=0.1111
OK surgeston => surgestone with dist=0.1111
OK eméthotrexate => methotrexate with dist=0.1538
OK luteny => lutenyl with dist=0.1667
OK témestat => temesta with dist=0.2500
OK purinéthol => purinethol with dist=0.1000
OK picroval => microval with dist=0.1250
KO esoprex => eprex with dist=0.2857
OK pantoprazol => pantoprazole with dist=0.0909
OK leelou => leeloo with dist=0.1667
OK allergenes => stallergenes with dist=0.2000
OK mnidril => minidril with dist=0.1429
OK diazepam => diazepam with dist=0.1250
OK metformin => metformine with dist=0.1111
OK monazole => monazol with dist=0.1250
OK gynergène => gynergene with dist=0.1111
```

For the active ingredient, only 25 fixes have been detected.

```
17 folliculum, folliculo
18 semaglutide, maglutide
19 manosonique, monosodique
20 nitr, nite
21 tranéxamique, tranexamique
22 rosavastatine, rosuvastatine
23 sultopril, sultopride
24 alrs, ars
25 etamsylate, tamsylate
```

The general vocabulary fixes up 430 words with relative high error rate (~25%): as accent encoding is badly handled by pysspellchecker module, I fixed it manually afterwards.

```
OK qui'l => quil with dist=0.2
OK siagnements => saignements with dist=0.18181818181818182
OK extremment => extremement with dist=0.1
KO ethpo => ethan with dist=0.4
OK disgestions => digestion with dist=0.18181818181818182
OK utilliser => utiliser with dist=0.11111111111111111
OK boulimies => boulimie with dist=0.11111111111111111
KO douelurs => douleurs with dist=0.25
OK hallucinogène => hallucinogène with dist=0.15384615384615385
OK cocceluche => coqueluche with dist=0.2
OK flushs => flashes with dist=0.16666666666666666
OK puisje => puisse with dist=0.16666666666666666
OK aujourdui => aujourd'hui with dist=0.11111111111111111
OK douloureuze => douloureuse with dist=0.09090909090909091
OK osteoporose => ostéoporose with dist=0.18181818181818182
KO govital => hopital with dist=0.2857142857142857
OK normalemnt => normalement with dist=0.1
```

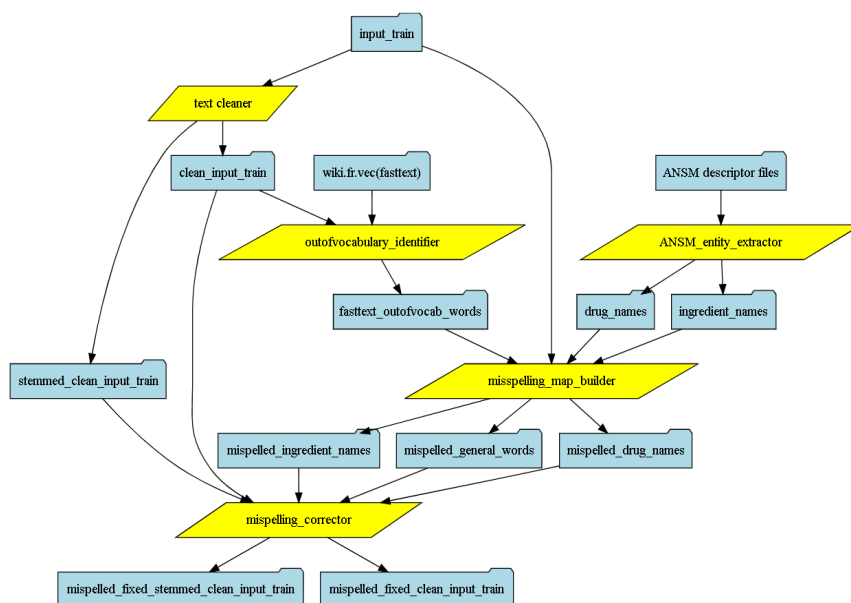
At the very end, it remains 583 unfixed words over an initial 2108 unknown words (25%): that corresponds to hard cases where the word is unexpectedly a concatenation of multiple word or transcribed phonetically.

```

539 pullile
540 spraypax
541 esteceque
542 jèmeré
543 anxios
544 'ai
545 présentent
546 enchainais
547 efezial
548 oedeme
549 acneique
550 depersonalisation
551 calmosine
552 demergent

```

Below diagram shows the different python notebooks (parallelogram in yellow) necessary to fix word misspelling: the blue folder represents the file consumed or produced by the python processing unit.



Lexical and Grammar Tagging

NER is a NLP process to tag text token with predefined categories (location, person, quantity, ...), permitting to count such entities as explanatory feature. Typically, distinct drug product counting may be a discriminating feature to predict the “drug interaction” label.

In fact, the Chinese NORP market has the three CARDINAL most influential names of the retail and tech space – Alibaba GPE, Baidu ORG, and Tencent PERSON (collectively touted as BAT ORG), and is betting big in the global AI GPE in retail industry space. The three CARDINAL giants which are claimed to have a cut-throat competition with the U.S. GPE (in terms of

POS (Part Of Speech) tagging is a process to markup text tokens with lexical categories (noun, adjective, verb, ...), enabling to compute tag frequency distribution as explanatory feature.

The stemming/lemmatization preprocessing is counterproductive to word embedding model learnt from corpora which haven't been stemmed or lemmatized upfront: they are so mutually incompatible and for DL scenario, I made use of word embedding excluding de facto this root normalization.

Classical Technique

Abstract

HMM (Hidden Markov Model) is a probabilistic and transitional graph modeling which is appropriate to model sequence of words. For example, it can learn on text corpus and predict POS tags but some research studies indicated that HMM is also applicable to text categorization with good performance: unfortunately, up to date HMM python implementation is missing.

The arguable fallback is to switch to statistical inference method like SVM, decision tree and so on, with the crucial loss of sequence awareness. To compensate slightly such discarding, the feature extraction/enrichment should include some handmade tricks trying to grasp some contextual information from the word sequence.

Feature Enrichment

This step adds a-priori extra features which may discriminate the label much more than the original features: they can be calculated from the text or can originate from external sources.

I incorporated above basic statistics giving insights on the text structure and composition:

- count of sentences
- count of words
- distinct count of drug name entities
- distinct count of active ingredient entities
- count of question marks (typically to identify specifically multi-intent label)
- individual count of interrogative pronoun entities (one column per pronoun: quand, qui, quoi, ou, comment, pourquoi, combien, quel(s)|le,...)
- distinct count of time entities (eg: jours, après midi, soir, année, 12h, mardi, samedi, temps....)
- distinct count of quantity entities (eg: 5mg, 10ml, ...)
- count of association entities (eg: et, avec, ou, ...)
- distance between interrogative pronoun and drug name entities
- distance between active ingredient and drug name entities
- distance between quantity and drug name entities
- distance between time and drug name entities
- distance between question marks and drug name entities

They are either **count-based or distance-based statistics**: distance variant is intended to catch the word context by measuring the relative distance between key entities. This computation needs to put in place the domain-based (list of distinct values) or custom regular expression NER (Named Entity Recognition) so that it's possible to locate the key entities in consideration.

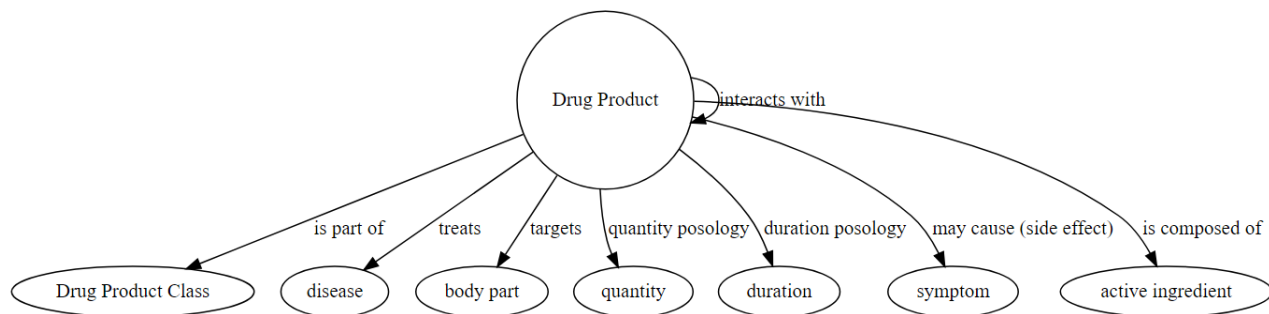
An extra calculated column is added to the train data frame per statistics as shown below:

```
XTrain['drugCount'] = XTrain['question'].map(lambda text: getEntityCount(text, drugNames, True))
XTrain['ingredientCount'] = XTrain['question'].map(lambda text: getEntityCount(text, ingredientNames, True))
```

```
XTrain['timeCount'] = XTrain['question'].map(lambda text : len(re.findall(timeRegex, text)))
XTrain['quantitiesCount'] = XTrain['question'].map(lambda text : len(re.findall(quantityRegex, text)))
XTrain['questionMarkCount'] = XTrain['question'].map(lambda text : len(re.findall(questionMarkRegex, text)))
XTrain['sentenceCount'] = XTrain['question'].map(lambda text : 1 + len(re.findall(sentenceSeparatorRegex, text)))
XTrain['wordCount'] = XTrain['question'].map(lambda text : getWordCount(text))
```

If the text sample has well identified drug product entities, it's valuable to extend the primary feature vector with relevant information related to these drug products.

I represent herein a dedicated **knowledge sub graph** centered around the drug product entity with some interesting relationships to other entities (quantity, human body part , ...).



Indeed, such related entities characterize well the drug product and they can improve the detection of the commonality between texts sharing same label: for instance, a drug product class (eg: antidepressant family) may raise particular questions.

Unfortunately, this knowledge graph model is not available publicly and should be built by our own: the ANSM provides online the full description of the drug usage indication in HTML format. Such resource can feed a learning system to extract above explanatory related entities.

I didn't implement this information extraction from ANSM source because it's a colossal workload which is incompatible with the project scope.

Last but not the least, I completed the feature enrichment with the **likelihood estimate that a text is associated a topic extracted by the NMF** algorithm.

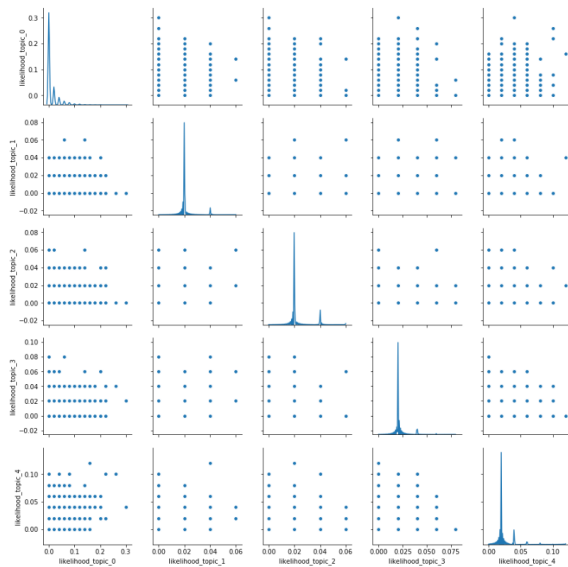
This probability is merely the ratio of matching word count between the text and the topic components over total number of topic components (I fixed empirically to 50).

Here's an overview of the extra columns representing the proportion of topic components used in the text:

```
topicLikelihoodFrame.head()
```

	likelihood_topic_0	likelihood_topic_1	likelihood_topic_2	likelihood_topic_3	likelihood_topic_4	likelihood_topic_5	likelihood_topic_6	li
0	0.08	0.0	0.00	0.02	0.02	0.02	0.06	
1	0.00	0.0	0.00	0.00	0.00	0.00	0.00	
2	0.02	0.0	0.00	0.00	0.00	0.00	0.00	
3	0.00	0.0	0.02	0.00	0.02	0.00	0.00	

I built a correlogram on the first 5 extracted topic likelihoods and I noticed no remarkable linear correlation: it indicates that topics are orthogonal/independent.



Feature Representation

The document (ordered set of sentences) should be converted into numerical vector because most of ML classifiers can only cope with numerical values and they don't care about symbol and semantic conveyed by the word.

First basic solution is the **BOW** (Bag Of Words) representation where each word of the vocabulary is defined in column and the text in row: the cell value stores the word frequency.

I didn't consider **n-gram** document representation because as mentioned earlier, the classical technique scenario doesn't employ sequence modeling like HMM which is able to treat n-gram structure.

The problem of the **BOW** representation is that rare term which in general discriminates well the document are under estimated in regards with commonly used but irrelevant terms (eg: generic verb, ...).

TF-IDF (Term Frequency Inverted Document Frequency) overcomes this pitfall by overweighting terms which are identified as rare for a given corpus.

The shortcoming is that such vectorization generates a very high dimensional space depending on the vocabulary size. We fall into the well-known **curse of dimensionality** where data distribution is extremely sparse making classification task inefficient when training size is too short.

The space dimension should be reduced consequently:

stop words and stemming processes already reduce upfront the vocabulary size

I applied the **PCA** (Principal Component Analysis) linear dimension reduction which keeps the top eigen vectors capturing the maximum of the data distribution variance: PCA is a process which is totally semantic unaware in contrary to word embedding I will tackle later on.

TF-IDF vectorization and PCA reduction produce a low dimensional numerical vector per document as below:

	0	1	2	3	4	5	6	7	8	9
0	0.160047	-0.080491	-0.043188	-0.011275	0.024025	-0.017478	-0.038179	0.056354	0.029396	0.020276
1	0.021063	0.161065	0.157428	-0.008642	-0.059919	0.053919	-0.057782	0.035386	-0.022574	0.113641
2	0.009776	-0.030112	-0.005735	0.011957	-0.008956	-0.011416	-0.042868	-0.094769	0.027011	-0.000518

Classification Modeling

Architecture

The classifier takes as input a feature space combining the reduced BOW representation and the handcrafted statistics:

	0	1	2	3	4	5	6	7	8	9	...	quandCount	quoiCount	commentCount	avec
0	0.160047	-0.080491	-0.043188	-0.011275	0.024025	-0.017478	-0.038179	0.056354	0.029396	0.020276	...	0	0	0	0
1	0.021063	0.161065	0.157428	-0.008642	-0.059919	0.053919	-0.057782	0.035386	-0.022574	0.113641	...	0	0	0	0
2	0.009776	-0.030112	-0.005735	0.011957	-0.008956	-0.011416	-0.042868	-0.094769	0.027011	-0.000518	...	0	0	0	0

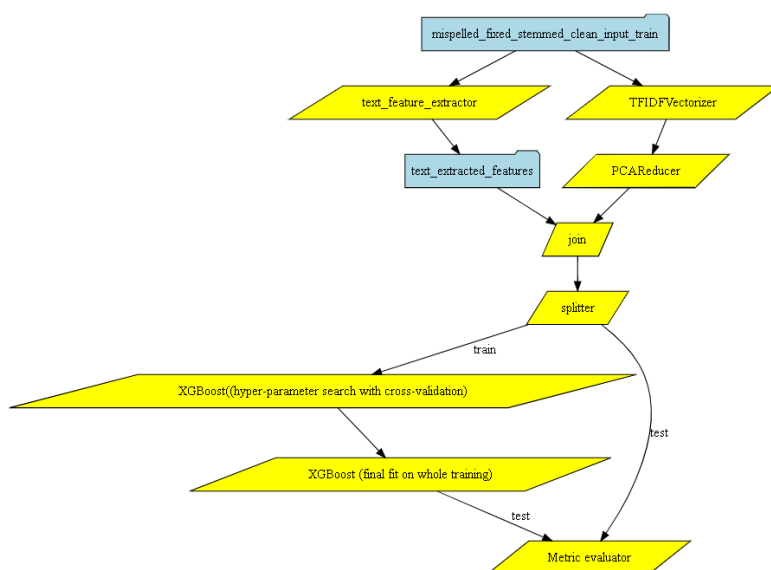
I bet on the **XGBoost** classifier delivering excellent accuracy in a reasonable time (it's multi-thread friendly): XGBoost is based on boosting ensemble technique combining sequentially weak classifiers (in general decision tree) where at each iteration, the weighting on incorreced classified observations is increased to enforce the next classifier to focus its attention on feature sub space with high error.

XGBoost comes up with many hyperparameters to tune: an inappropriate selection usually leads to suboptimal model.

I followed the standard methodology and best practices:

- find out the optimal hyper-parameters by testing different combinations. I retained the one delivering the best accuracy on unseen dataset (validation) with **cross validation** enable as training is very small
- fit the final model with the above fixed hyper-parameters on the whole training and assess the generalization error on test

Here's the learning pipeline for the classical technique track:



Hyperparameter Search

I concentrated my attention on the following parameters which are the most instrumental to the final accuracy:

max_depth

This parameter drives the decision tree complexity to partition the feature space.

A low value usually prevents from overfitting and favor the weak learner synergy.

I tested empirically 3 values: 4 , 6 and 8.

min_child_weight

Under the threshold, the learner stops splitting and generates a leaf node.

It controls as well the tree complexity and consequently the overfitting.

I tested empirically 3 values: 2, 5 and 10.

n_estimators

This parameter sets the maximum number of stacked trees.

I fixed it empirically to 100.

learning_rate (eta)

It controls an important parameter of the gradient descent optimizer. Too small value leads to extremely slow optimization whereas too high value would miss global extrema.

I tested empirically 2 values: 0.05 and 0.1

cross validation fold

Cross validation ensures a more reliable generalization error indicator which is not biased by a particular split (test set). It's valuable typically in imbalanced label or small dataset situation: so it's applicable to this data challenge.

I fixed it empirically to 4: high value increases linearly the search time.

Other parameters settings rely on the XGBoost defaulting to avoid excessive processing time caused by the grid search combinatory explosion: by crossing max_depth, min_child_weight, learner rate and cross validation fold, it represents 72 (3x3x4x2) learning units to reveal the optimal parameter values.

For the final model fit, I set up the early stopping parameter to 10 in order to avoid useless extra tree stacking.

Result Analysis

The accuracy on test is low:

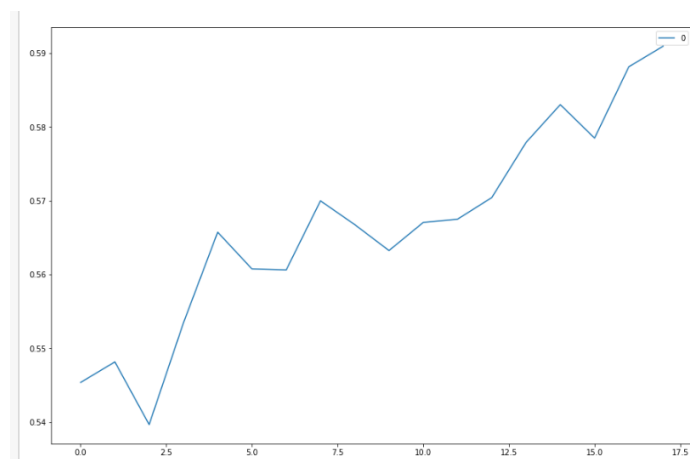
	micro F1-score	macro F1-score	support
avg / total	0.63	0.44	1205

macro F1-score ignores label imbalance and corresponds to the average of per label F1-scores, whereas micro F1-score is evaluated from the whole confusion matrix with no intermediate F1-score evaluation per label.

In the context of the challenge, micro F1 score is a more relevant metric as the label is not balanced.

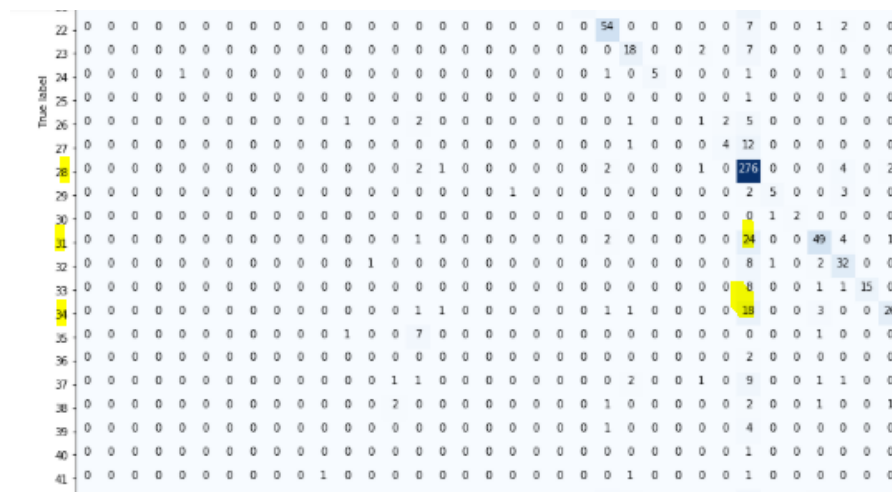
This figure displays the score for each hyperparameter selection and confirms that tuning is a key factor of accuracy (score ranges from 0.54 to 0.59): the best score is obtained with learning_rate=0.1, max_depth=8 and min_child_weight=10.

The grid search with CV (fitting 4 folds for each of 18 candidates, totaling 72 fits) took 2 hours to find the optimal parameters.



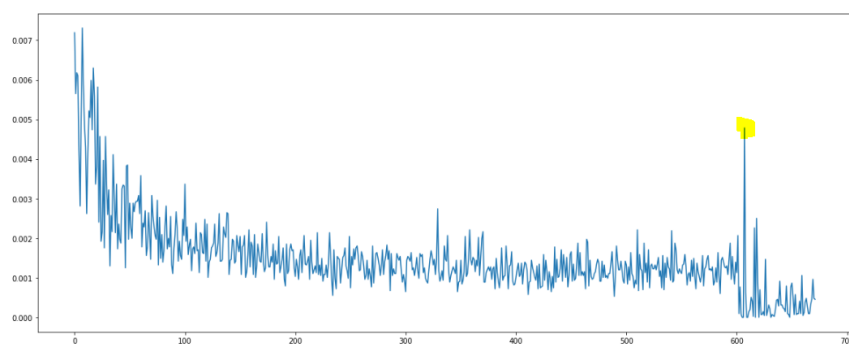
When analyzing the 51x51 confusion matrix, the highest confusion occurs on intent 31 (drug>disease **indication**) which is incorrectly predicted to intent 28 (drug>disease **contraindication**).

The high confusion error between 31 and 28 confirms the weakness of BOW representation I used here: 2 entities occur simultaneously (drug product and disease) but the context linking them is not understood (indication vs contraindication).



The visualization of the variable importance shows up that:

- the variable importance is correlated with the variance level of the PCA components (components capturing the most variability are ranked first in the feature importance list)
- the manually defined statistical features (at the tail) are poorly explanatory except an outlier highlighted in yellow: unexpectedly, wordcount is relatively important!



```
9]: namedVarImportances = pd.DataFrame({'feature_names': namedVarImportances[600:]})
```

	feature_names	importance
600	ID	0.001130
601	drugCount	0.002068
602	ingredientCount	0.000096
603	timeCount	0.000770
604	quantitiesCount	0.000072
605	questionMarkCount	0.000000
606	sentenceCount	0.000000
607	wordCount	0.004785
608	combienCount	0.000625
609	pourquoiCount	0.000000

Deep Learning Technique

Abstract

DL is widely recognized as an universal estimator capable of fulfilling any sort of learning task from feature representation to the predictive modeling within a single neural network. The key strength of this all-in-one learning is that the loss optimization to find out the best modeling parameters (weights, ...) operates consistently across all functional layers regardless of their respective purpose (embedding, decision making, ...). In contrast, with traditional method, feature representation and classification are 2 sub tasks which are engineered/optimized separately.

I specifically looked at its sequence modeling capacity carried by 2 architecture types:

- **RNN** (Recurrent Neural Network with **LSTM** (Long Short Term Memory) unit
- **CNN** (Convolutional Neural Network)

The hybrid option mixing up CNN and RNN is not considered here for simplicity sake even if some practitioners recommends this winning combination to get cutting edge performance.

Furthermore, DL also provides a very good support of word embedding which can be combined nicely with above architectures as upstream layer.

Feature Enrichment

I intentionally excluded extra features to verify how a DL sequence modeling can give some good results without manual contributions (statistics on text, ...).

Feature Representation

Sequential representation

As the predictive modeling layer is sequence aware, the text representation should be **n-gram** where n is the number of words to keep: if the number of words from document is insufficient, it's necessary to apply a zero padding to get at the end a fixed sequence length for all documents.

As mentioned in the "Data Exploration" section, lengthy document usually starts with the description of the question context and ends up with concrete question (eg: "Je suis suivi par un médecin ... Qu'est ce que c'est recommandé?"). It would make sense as the document can be truncated due to the fixed sequence length constraint to keep the n-th last words and not the n-th first words to not lose the question part.

In short, each document is shaped as a `fixed_sequence_length x vocabulary_size` matrix: again, vocabulary size can be huge leading to inappropriate high dimensional feature space and a dimension reduction is mandatory.

Word embedding

Text corpus

Instead of applying a generic PCA, a better alternative is the popular **word embedding**: it's an unsupervised method which learns on a very large text corpora to optimize a lower dimensional vector representation where words sharing similar context (within a sentence) are close to each other. The wonder of this dimension reduction is that vector proximity is governed by semantic similarity.

Word embedding is implemented in a DL flavor (Word2Vec or FastText) and in a non-DL way too (GloVe project) with nearly similar performance.

The question now is to choose the **text corpus** used to build this embedding model:

- consume directly model tediously pre-trained by the GAFA companies. Such model is based on very large general vocabulary but probably miss domain specific vocabulary (our study case in fact)
- build a custom embedding from the POSOS challenge corpus. It overcomes the domain specific vocabulary lack (drug product name, ...) but it's not complete and robust enough considering the small training dataset with many misspelling/incorrectness in the text.
- perform a model transfer from GAFA base with specific vocabulary coming from POSOS corpus. In practice, this ideal solution is undoable because it's required tremendous amount of GPUs to rebuild a merge embedding model combining general and specific vocabulary.

I finally experimented the custom and general embedding models and discarded transfer model option. For the general vocabulary embedding, I opted for the 300-dimensional **FastText** model which is gracefully available in French language.

In conclusion, the training dataset is represented as a $n \times k \times v$ numerical matrix where:

- n is the number of observations (document)
- k embedded space dimension
- v fixed sequence length

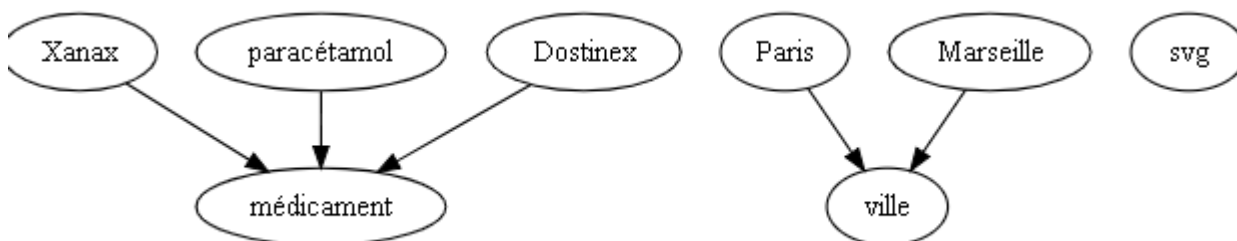
OOV handling strategy

Embedding layer can only deal with word which exists in its vocabulary: if the learnt corpus and the corpus to vectorize are dissimilar, OOV is potentially frequent. The usual solution is to encode unknown words into a random embedded vector at the risk of generating noisy feature representation.

A more elegant alternative is to merely project such unknown words into its **hypernym** entity (having a type-of relationship with the concerned word) guaranteeing a semantic proximity in the embedded space for entities of the same class/hypernym:

- all drug product entities (eg: Xanax, Abboticine) is replaced by 'médicament'
- all active ingredients (eg: Acabavir) is replaced by 'médicament'

To not completely lose the subtle distinction between entities sharing the same class, I added a very small stochastic variation vector based on the entity name so that all 'Xanax' entities have exactly the same vector and are also close to 'Paracétamol' entities.



The custom extension of FastText model is managed by the `fasttext_embedding_extension_builder.ipynb` script.

DNN Architecture

I setup a test with DNN (Dense Neural Network) which is not a sequence learner, as a comparison baseline for more sophisticated architectures like RNN and CNN. It would give a good hint on the relative performance gain with sequence awareness in the modeling procedure.

Moreover, the embedding layer is built from the POSOS corpus to define again a comparison baseline to measure the relative gain (or loss) when opting for general purpose corpus.

The concrete architecture is described by the summary output generated by Keras:

Layer (type)	Output Shape	Param #
embedding_2 (Embedding)	(None, 30, 300)	3000000
flatten_1 (Flatten)	(None, 9000)	0
dense_1 (Dense)	(None, 200)	1800200
dropout_1 (Dropout)	(None, 200)	0
dense_2 (Dense)	(None, 51)	10251
Total params: 4,810,451		
Trainable params: 4,810,451		
Non-trainable params: 0		

There are 2 dense layers with different activation functions: relu at the first layer and softmax at the decision layer. The dropout layer is placed between to introduce some random perturbation to combat overfitting. I set the embedding dimension to 300 in accordance with the pretrained FastText model.

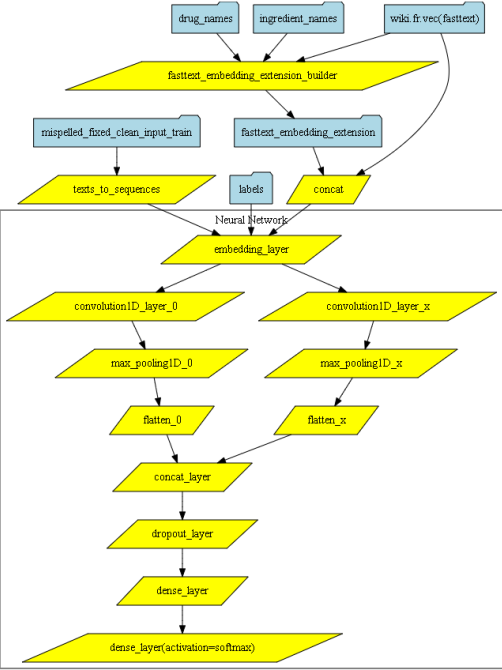
CNN Architecture

In **Computer Vision (CV)**, convolution operation is well known to be remarkable in extracting the high level representation of an image by applying a sliding window filter and computing consequently an average value for each filter position as output. These convoluted values are then activated with usual non-linear function and down-sampled thanks to the pooling layer. This pixel-wise processing is inspired by how the visual cortex analyzes the signal sent by the eye receptors.

Surprisingly, such biological inspiration also works well to catch the structural sense of word sequence in NLP. The convolution operates in a 1-dimensional array (word sequence) instead of 2D (pixel matrix) in CV. The sequential filter enforces the neural network to focus its attention on local context which establishes connection between words.

Even if some research studies demonstrate that convolution is expressive enough to cover embedding contribution, I setup my CNN architecture with embedding layer upfront. As usual, some dropout layers are intermittently inserted into the neural network.

The paper **REF1** (see appendix) recommends building many **parallel convolutional layers** with different filter sizes and/or strides whose outputs are then concatenated to each other and this is the resulting predictive pipeline:



Layer (type)	Output Shape	Param #
input_6 (InputLayer)	(None, 30)	0
embedding_6 (Embedding)	(None, 30, 300)	2429400
reshape_6 (Reshape)	(None, 30, 300, 1)	0
dropout_11 (Dropout)	(None, 30, 300, 1)	0
conv2d_17 (Conv2D)	(None, 28, 1, 100)	90100
conv2d_18 (Conv2D)	(None, 27, 1, 100)	120100
conv2d_19 (Conv2D)	(None, 26, 1, 100)	150100
max_pooling2d_17 (MaxPooling2D)	(None, 1, 1, 100)	0
max_pooling2d_18 (MaxPooling2D)	(None, 1, 1, 100)	0
max_pooling2d_19 (MaxPooling2D)	(None, 1, 1, 100)	0
flatten_17 (Flatten)	(None, 100)	0
flatten_18 (Flatten)	(None, 100)	0
flatten_19 (Flatten)	(None, 100)	0
concatenate_6 (Concatenate)	(None, 300)	0
dropout_12 (Dropout)	(None, 300)	0
dense_11 (Dense)	(None, 100)	30100
dense_12 (Dense)	(None, 51)	5151

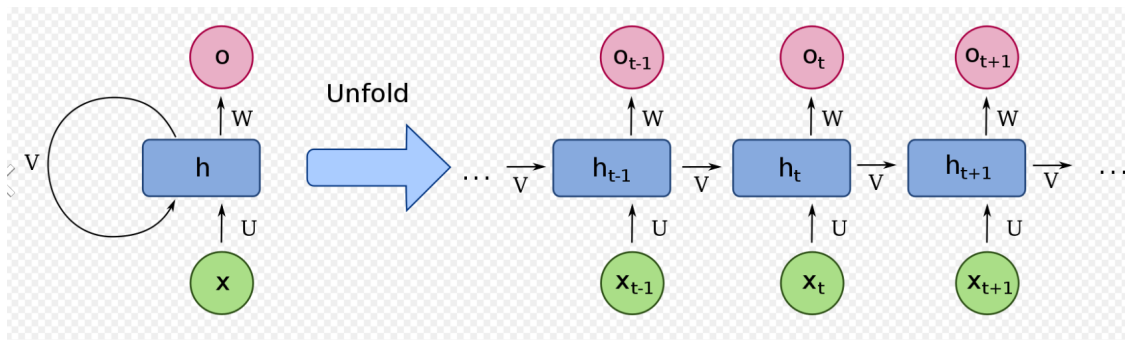
Another CNN architecture alternative proposed by REF2 (see appendix) is to define a **sequential layout** of the different CNN layers as illustrated below:

Layer (type)	Output Shape
embedding_3 (Embedding)	(None, 15, 300)
dropout_5 (Dropout)	(None, 15, 300)
conv1d_7 (Conv1D)	(None, 100, 298)
max_pooling1d_7 (MaxPooling1D)	(None, 100, 298)
conv1d_8 (Conv1D)	(None, 100, 295)
max_pooling1d_8 (MaxPooling1D)	(None, 100, 295)
conv1d_9 (Conv1D)	(None, 100, 291)
max_pooling1d_9 (MaxPooling1D)	(None, 50, 291)
flatten_3 (Flatten)	(None, 14550)
dropout_6 (Dropout)	(None, 14550)
dense_5 (Dense)	(None, 100)
dense_6 (Dense)	(None, 51)

RNN Architecture

Recurrent Neural Network architecture tries to leverage this sequential information from the sentence with a special layout where each layer at the i-th position is fed with the i-th element of the input sequence and the output of the direct preceding layer (i-1 th position): each layer captures somehow the hidden state (memory) of the preceding sub sequence of inputs (words).

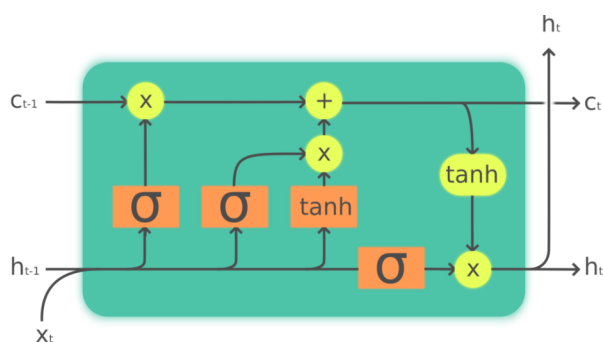
RNN can be **bi-directional** (instead of forward only) where the i-th layer also depends on the computational output of the direct successor (i+1 th position).



(source: Wikipedia)

In practice, the simple layer computational unit (“h” blue box in above diagram) exhibits inability to catch long term dependencies to distant input $x(t-n)$ where n is significant high, due to the vanishing gradient problem.

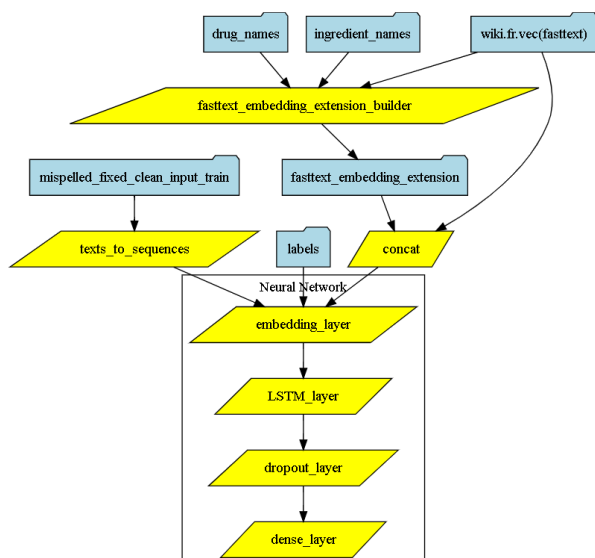
The **LSTM** (Long Short Term Memory) cell unit has been invented to treat this long-tailed sequential dependency we can find in multi-sentence text analysis where the context may be specified upfront far away from the concerned word.



(source: Wikipedia)

This processing unit adds a secondary flow (upper stream) to update gradually the memory (cell state denoted as $C(t)$) with contributions controlled by several input gates (shown at the bottom).

The DL architecture with the word embedding extension is represented below:



Layer (type)	Output Shape	Param #
embedding_6 (Embedding)	(None, 30, 300)	2430000
bidirectional_2 (Bidirectional)	(None, 600)	1444800
dense_9 (Dense)	(None, 40)	24040
dropout_5 (Dropout)	(None, 40)	0
dense_10 (Dense)	(None, 51)	2091
Total params: 3,900,931		
Trainable params: 1,470,931		
Non-trainable params: 2,430,000		

CuDNNLSTM is a very convenient Tensorflow/Keras class automating the construction of the recursive processing unit pattern: furthermore, it's based on the CuDNN Nvidia library boosting learning time by a factor of 10.

The bidirectional mode has been enabled just by wrapping the CuDNNLSTM construct with Bidirectional.

```

trainable=False))

# recurrent network layer
model_lstm.add(Bidirectional(CuDNNLSTM(embedding_out_dims)))

```

Hyperparameter search and Architecture Sizing

DNN

No particular tuning for DNN baseline: nevertheless, when augmenting number of dense layers, the final accuracy is not improved at all and at the end, I kept going with 2 dense layers whose one of them is aimed to produce the classification probability.

Sequence Max Length

When encoding the document into matrix, a sequence max length should be defined: if the document is shorter than this max length, the sequence padding is applied.

I tested several values (15, 20, 25 and 30) considering the word count mean (9) and the high variance (13) and kept the one producing the best score with a LSTM architecture: I observed unexpectedly that high value (30) in some architectures performs well.

Parallel CNN

Based on the empirical recommendation from papers **REF1** and **REF3**, filter sizes should be 3, 4 and 5 with a max pooling size set to 1. Even if REF3 asserts that adding dropout brings marginal gain, it places intermittently some regularization layers. I played with different values for filter size and max pooling size and it doesn't impact positively the final accuracy.

Sequential CNN

I started with parameter setting found from reference **REF2** but the result was not compelling: filter sizes = (5,5,5) and max pooling sizes = (5,5,35).

After several manual attempts, (3,3,3) and (3,3,15) values appear to deliver better performance.

RNN / LSTM

LSTM has quite few numerical hyperparameters and I definitely trust the default setting. The unique tuning I did is to compare the single and bidirectional RNN: it turns out that there was no major performance difference and I used finally the single one which is slightly faster.

DL Implementation and Execution

I have implemented all DL networks with the high level **Keras** model which runs on top of **Tensorflow**. Keras provides a very friendly API that simplifies dramatically the neural network construction by hiding all the technical boiler plates (TensorFlow session handling, many default parameters are set, ...).

Here's an example of Keras code where the layers are built in a very concise manner thanks to wrapper objects like Convolution2D, Activation and so forth.

```
# build neural network
model_lstm = Sequential()

# dimension reduction layer
model_lstm.add(
    Embedding(
        len(tokenizer.word_index)+1,
        embedding_out_dims,
        weights=[embedding_matrix],
        input_length=sequence_length,
        trainable=False))

# recurrent network layer
model_lstm.add(Bidirectional(CuDNNLSTM(embedding_out_dims)))

# classification hidden layer
model_lstm.add(Dense(hidden_dims, activation="relu"))

# random node inactivation
model_lstm.add(Dropout(dropout_ratio))

# normalization layer
model_lstm.add(Dense(num_classes, activation='softmax'))

model_lstm.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])

model_lstm.summary()
```

Nevertheless, it's still possible to access to the specific APIs of the underlying concrete framework to configure the execution parameters as shown below.

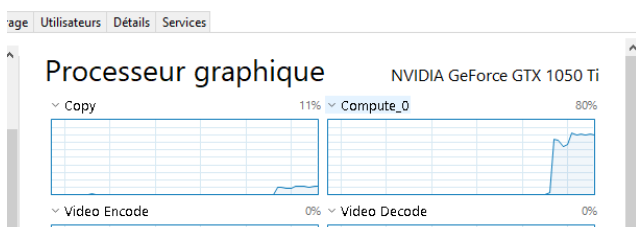
```
# tensorflow technical setting
config = tf.ConfigProto(device_count={"CPU": 32})
gpu_options = tf.GPUOptions(per_process_gpu_memory_fraction=0.8, allow_growth = True)
config=tf.ConfigProto(gpu_options=gpu_options,allow_soft_placement=True)
keras.backend.tensorflow_backend.set_session(tf.Session(config=config))
```

With CUDA/GPU activation, Tensorflow runs really faster than using the regular CPU (even with a powerful 32-core machine) but it asks for a particular attention on the DL configuration: when inappropriately parameterized, it causes in worst case crashes (allow_growth parameter should be defined) or memory allocated error.

```
E tensorflow/stream_executor/cuda/cuda_blas.cc:459] failed to create cublas handle: CUBLAS_STATUS_ALLOC_FAILED
E tensorflow/stream_executor/cuda/cuda_blas.cc:459] failed to create cublas handle: CUBLAS_STATUS_ALLOC_FAILED
E tensorflow/stream_executor/cuda/cuda_blas.cc:459] failed to create cublas handle: CUBLAS_STATUS_ALLOC_FAILED
E tensorflow/stream_executor/cuda/cuda_blas.cc:459] failed to create cublas handle: CUBLAS_STATUS_ALLOC_FAILED
W tensorflow/stream_executor/stream.cc:2818] attempting to perform BLAS operation using StreamExecutor without BLAS support
```

GPU memory consumption is sensitive to the training size and the batch size parameter: higher value means higher GPU memory but batch size selection impacts the optimizer behavior and consequently the resulting model.

My GTX1050Ti graphic card with 4Gb memory is a bit short to handle mid-complex DL training: setting the memory threshold to 80% implies more data movement (Copy operation) between the motherboard and GPU memories.



Result Analysis

Here's the overall classification score per architecture:

	Micro F1-score	Macro F1-score
DNN / Custom embedding	60%	40%
Parallel CNN / Fasttext embedding	67%	48%
Sequential CNN / Fasttext embedding	61%	45,7%
RNN-LSTM / Fasttext embedding	65%	44%

Such results have been achieved by testing manually few combinations of hyperparameters, in opposition to the systemic parameter grid search for XGBoost estimator.

LSTM and parallel CNN outperform slightly sequentially CNN and DNN architectures: sequential CNN exhibits unexpectedly very disappointing score.

The CNN parameters such as filter size, stride and pool size are difficult and non-intuitive to define optimally and the resulting performance is pretty sensitive to these hyper parameters.

LSTM learning curve indicates that the model is overfitting on train beyond 12 epochs but there's no observed accuracy improvement on validation set (the accuracy chart is stationarized around 63% beyond 6 epochs).

This means that this LSTM (in general any neural networks) is capable to fit perfectly any data samples but its learning is badly applicable to unseen dataset (validation): this is typically an overfitting risk.

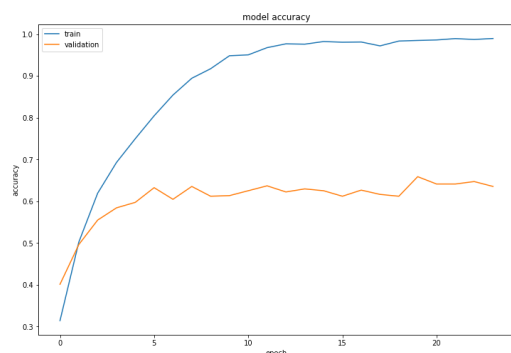
For all architectural candidates, I enabled the early stopping mechanism on validation loss as a call back when fitting the model.

```
# stop criterion to avoid overfitting
call_back_early_stopping = keras.callbacks.EarlyStopping(
    monitor='val_loss',
    min_delta=0,
    patience=patience,
    verbose=0,
    mode='auto',
    baseline=None)

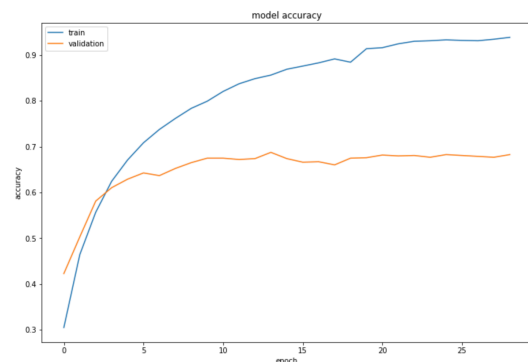
# learn !
model_lstm.fit(
    XEncodedTrain,
    np.array(YOneHotEncodedTrain),
    validation_split=0.10,
    epochs=num_epochs,
    verbose=2,
    callbacks = [call_back_early_stopping, call_back_board])
```

CNN and LSTM architectures need more iterations than the simple DNN to reach good fit on train. The learning curve shape is more eradic when reducing the batch size.

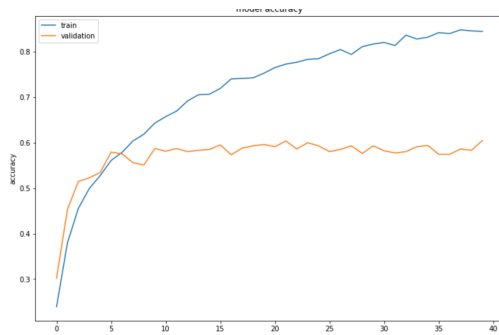
LSTM/RNN learning curve



Parallel CNN learning curve

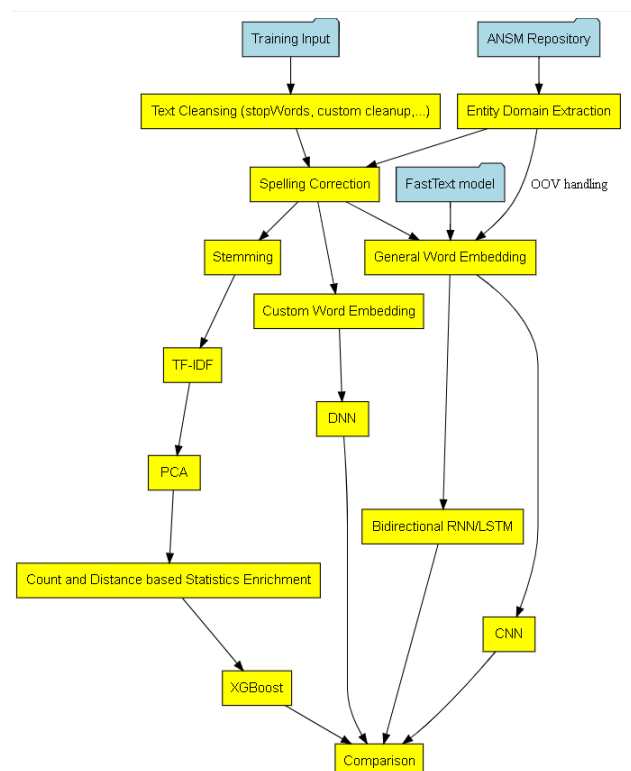


DNN learning curve



Comparative Study

I represent below the overall comparative experiment where 4 modeling candidates are evaluated holistically (accuracy, interpretability, ease of use, ...).



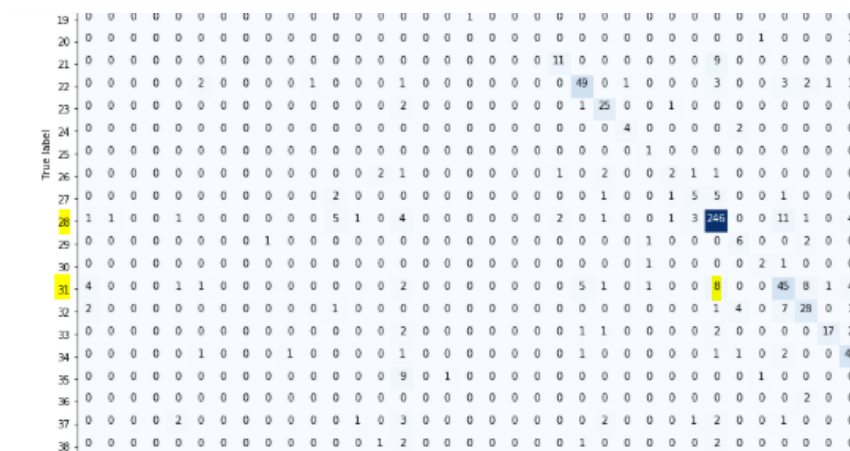
Model Accuracy

LSTM/parallel CNN architectures deliver a better performance (65-67% micro F1 score) than classical technique (63% micro F1 score).

Nevertheless, it represents only 2-4% gap which is not really significant: this disappointing difference may be explained by the empirical and non-expert choice I did on the hyperparameters and architecture of the neural network.

The scores on test I measured should be considered with caution: indeed, I hold out 15% of the training dataset (8000), representing only 1200 rows for 51 labels. It would be fairer to do a cross validation error measurement on test: I suspected that the models tend to overfit with significant standard deviation on the fold errors.

A noticeable point when looking at the confusion matrix, NN model does a better job than XGBoost to distinguish drug>disease indication and contraindication topics: the sequence/context awareness seems to pay off.



Model Interpretability

DL has the good reputation to provide excellent prediction accuracy when the architecture engineering is well conducted, but the theoretical/mathematical foundation is not rock solid yet (some mathematicians are still working on).

Additionally, model interpretability (feature importance, individual contribution at prediction time) is not supported natively with DL, except using exogeneous explainer like LIME (Local Interpretable Model-Agnostic Explanations) which can provide some model interpretations agnostically.

DL is very versatile and flexible to fulfill various modeling schemes but choosing the appropriate architecture or the optimal hyper parameters are often based on an empirical approach with few formal guidance: when a given setting produces a better accuracy, it's really tricky to find an explanation.

XGBoost is roughly the opposite: better native interpretability support (approximation of individual contributions, feature importance, ...), better math foundation but it's lacking on complex modeling support (sequence learning for instance). The fact that I defined manually extra features which are semantically explicit contributes as well to ease the model interpretability.

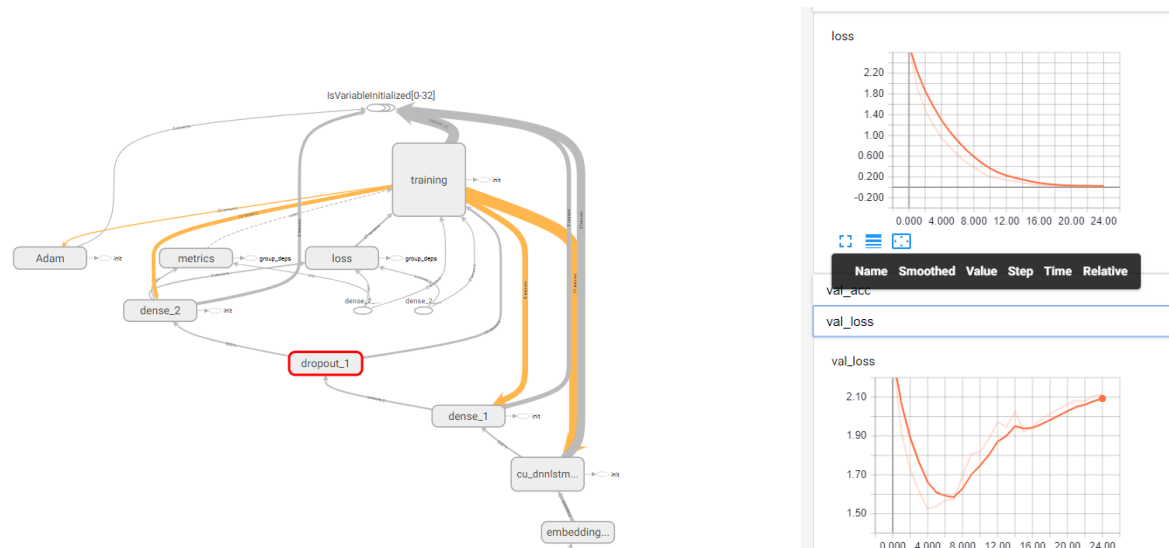
This is the well-known trade-off between model accuracy and model interpretability.

Tooling

From a practical standpoint, the classical technique with scikit-learn framework offers a very good level of tooling to implement rapidly common practices/methodologies like hyper parameter search and cross validation: XGBoost learner is plainly compliant with scikit-learn framework. Moreover, scikit-learn comes up with a complete and comprehensive set of features (text processing, LDA, NMF, feature extraction, many classifiers/regressors, ...) explaining its popularity.

On the other hand, Keras framework supports only Deep Learning paradigm and provides a very concise neural network construction API by favoring configuration by exception (default parameters). It hides the complexity of DL but it doesn't unfortunately solve it by automating and applying heuristics to assist on selecting the ton of hyper parameters behind the scene.

Nevertheless, when combined with Tensorflow, refining the architecture and parameters is greatly backed up by TensorBoard visualization which gives nice insights on the detailed neural network flow and on the different learning curves (it's a pity that it's not possible to draw both learning curves on train and validation in the same figure).



Sustainability

XGBoost/scikit-learn and Keras/Tensorflow combos are both very active open source projects: contributions to Keras/Tensorflow mainly come from Google organization whereas XGBoost/scikit-learn project is the fruit of academic field.

Nevertheless, DL technology is much more popular to tackle unstructured data (video, image, text, ...) and benefits a larger support of the ML researchers in the area of NLP.

Improvement Tracks

There are for sure a lot of improvement rooms on the learning procedure I have elaborated so far. Here are some possibilities to enhance the model accuracy, I feel like to implement if I had more times and means (GPU farm, ..).

Spelling Correction

Spelling correction relies on the Levenshtein distance and the fix suggestion is not influenced by more accurate criteria:

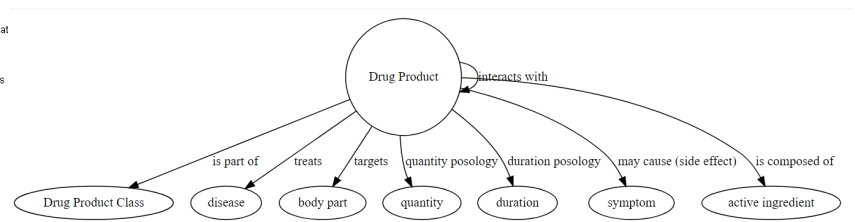
- in case of equality, favor probable words where the difference is located on the accent (eg: reveiller vs réveiller)
- take into account of the word frequency observed in the drug specific corpus (drug/medical) (eg: for misspelled someil, sommeil is more frequent than soleil in the drug question corpus)
- favor phonetically close words (eg: méson should be fixed into maison and not téton)

```
5 simplemen, simplement
7 dadaptation, adaptation
8 anxyolitique, anxiolytique
9 peremption, perception
10 allergie, allergie
11 someil, soleil
12 comprimés, comprises
13 ovulé, ovule
```

Named Entity Recognition

I have underexploited the ANSM repository providing in particular medication guide per drug. Information extraction can be performed quite easily by leverage the structure of the HTML page: for instance there's a dedicated/fixed section on adverse effect ("Quels sont les effets indésirables ..?") with a formatted list of undesired consequences.

The screenshot shows the ANSM website with a medication guide for XANAX 0.25 mg. The page is in French and includes sections on duration of treatment, frequency of use, and a list of potential side effects. The side effects listed are: dépression, sédation, somnolence, and difficulté à coordonner certains mouvements.



With a more complete knowledge graph on drug, it would be possible to build a wider NER system extending the basic one I actually built (only drug product and active ingredient entities) with drug product class, adverse effect, disease [contra]indication, ...

Another trick to get a more satisfactory French NER system is to use the English NER combined with an English to French translator.

Count/Distance based Statistics

With above extended NER system, it's worth to compute extra statistics on the new entities. For instance, count on adverse effect entities present in the sentence may be informative to explain the target.

Word Embedding

Custom embedding has been built on a too small corpus (training) and the result is badly robust. It would be valuable to extend this corpus with:

- test corpus (input_test.csv) even if there's no label (embedding is unsupervised)
- external source (eg: doctissimo.fr web site hosts discussion forum on drug)

Better solution is to merge above corpora with the ones used by FastText model and build our own embedding model: such learning involves a tremendous amount of resource (GPU, memory) and processing time.

OOV Handling

I proposed a "better-than-random" handling in case of out of vocabulary: project the drug product entities into its class/hypernym ("médicament") with a very small stochastic perturbation.

To reduce the taxonomical loss, a possible enhancement is to leverage the drug name class provided by the above extended NER system: the drug product entities would be converted into their respective drug class ("antidépresseur", ...).

Early Stopping With Grid Search CV

Scikit-learn GridSearchCV cannot leverage the early stopping of XGBoost on the fold (validation set) GridSearch wrapper defines internally: indeed, XGBoost.fit() asks for an explicit DMatrix for eval_set which is used by the early stopping mechanism.

This inability to use the early stopping on the fold during the grid search means that in a second time, once the best hyper-parameters are found, we need to run a XGBoost.fit() with early stopping on validation set in order to determine the best early stopping value.

```
gridSearch = GridSearchCV(
    estimator=gbm,
    fit_params = None,
    param_grid = grid_parameters,
    cv=4,
    verbose=1)

gridSearch.fit(mergedXTrain, YTrain.intention)
```

```
fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None,
    verbose=True, xgb_model=None, sample_weight_eval_set=None, callbacks=None)
```

Neural Network Tuning

For the classical method, Scikit-learn framework offers convenient wrapper to tune the hyper-parameters with cross-validation. On the DL side, Keras doesn't support natively cross validation nor kind of grid search wrapper to ease the execution and scoring of different hyperparameters combination.

Having said that, it's still feasible to write ad-hoc python in order to simulate the equivalent of grid search with cross validation: the practical pain point is that DL learning unit is slower than XGBoost.

Other Modeling Candidates

I purposely imposed that each candidate relies on an unique modeling principle to make the comparison more academic and distinctive. The state of art is to combine all the techniques in the hope to provoke a synergy effect where each estimator strength would overtake on average.

For example, associating RNN/LSTM and CNN in the same neural network would be a good candidate to be tested. For classical technique, building a general vocabulary embedding model with matrix factorization (similarly to GloVe) and combining it with XGBoost are likely worthwhile.

HMM (Hidden Markov Model) solution was also discarded unfairly whereas it's a sequence learning.

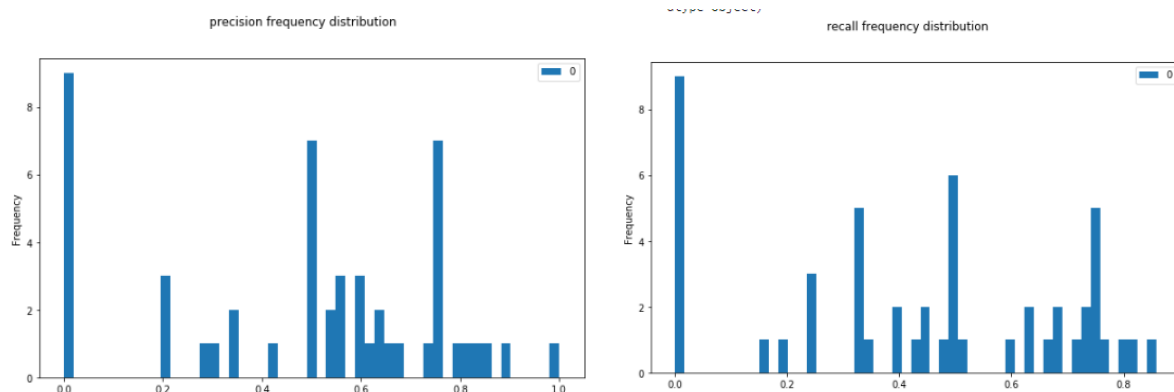
Too Few Samples and Too Many Target Effects

The training dataset has improper characteristics to get good learning level and cumulates 2 major failures:

- too few samples (~85% of 8000) which sounds insufficient for the data greediness of deep learning network where hundred thousands of parameters have to be determined.
- too many and imbalanced target labels (51 labels and 273 as frequency standard deviation)

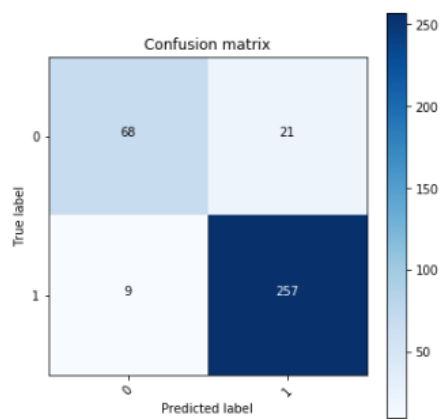
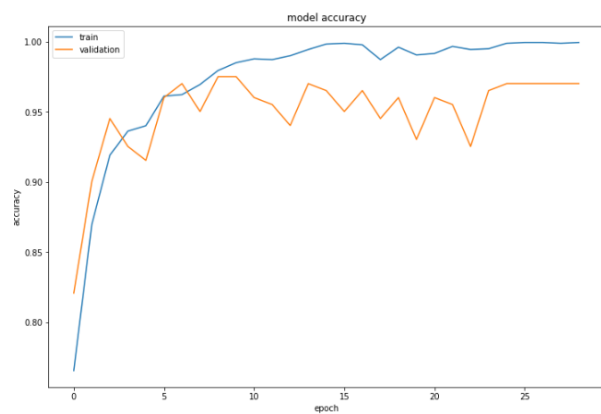
Moreover, I hold out ~ 15% of train for measure the generalization error on test: that represents 1205 samples as support, which are again not enough to compute a robust error on unseen data.

For instance, the LSTM-based learner produces consequently an imbalanced recall / precision on test.



I suspected that the optimizer used by XGBoost or TensorFlow tends to focus on the precision of dominant labels at the expense of minor labels: the important gap between the micro and macro F1 scores endorses this assumption.

When simplifying the multiclass problem into a **binary classification** by keeping the 2 most frequent labels, LSTM accuracy is improved dramatically with **92% micro F1-score!**



Conclusion

In this present paper, I describe a walkthrough return of practical experience on text categorization problem in a **DL (Deep Learning)** and non-DL fashion (a.k.a. classical method).

First of all, the POSOS problem is very challenging because the learning materials are not sufficient (8000 questions) for a multiclass classification task and the writing style is actually familiar with many misspellings. In addition, the high cardinality of target makes the problem harder: indeed, the good performance in text categorization exhibited in different research papers or blogs deals with binary classification (sentiment analysis) or reasonable cardinality in multiclass (less than 10 labels).

The classical method involving the famous XGBoost classifier delivers poor modeling performance (63% micro F1-score on test) for various specific reasons:

- no available word embedding model in French and a mere generic PCA as dimension reduction
- too limited count/distance bases statistics due to the lack of named entities (only drug and active ingredient)
 - this feature extraction is a fallback supposed to compensate the XGBoost inability to model sequence

DL scenario plays with several architectures (DNN, CNN, LSTM): unsurprisingly, as a native sequence modeling, RNN/LSTM overtakes slightly other variants with a 65% micro F1-score on test.

I intuitively hoped a bigger gap with classical method scoring at 63%: it's probably due to lack of hyper-parameter tuning or beginner's mistake on architecture choice.

On the other hand, traditional method with the combined Scikit-learn/XGBoost offers a better interpretability (feature importance, individual contributions) and methodology support (cross validation, grid search).

Even if the performance difference is not so important (4% micro F1-score at best), Deep Learning turns out to be the winning ML technique in NLP task solving: it's presently a very active research domain within the Data Science community as DL expressiveness/capacity are unbound similarly to the brain plasticity.

This short study with non-conclusive performance result, has at least the educational benefit to make me practice on a large spectrum of domains that a data scientist should master:

- data analysis to experimental result debriefing
- unsupervised (word embedding) vs supervised (classifier)
- training implementation and execution (GPU activation, python programming, AWS computing infrastructure, ...)
- Deep learning architecture variety (RNN, CNN, DNN, ...)
- general best practices on classification task (cross validation, early stopping, hyper-parameter search, ...).

Appendix

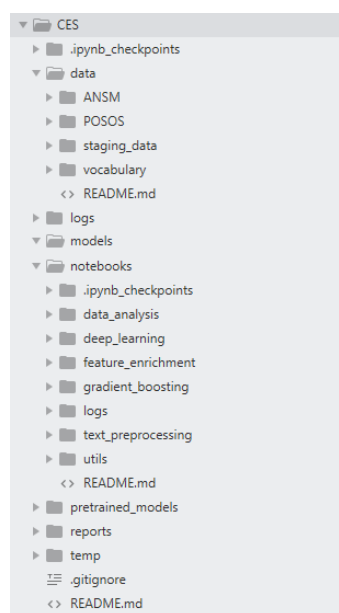
Github project

The project is available from the public github repository at the following URL:

<https://github.com/jhuu32/CES>

It contains Jupyter notebooks allowing to reproduce all learning experiments mentioned in this report: the only missing artifact is the FastText embedding model which is too large to be pushed into github (2Gb), but the README.md gives the necessary information to download it.

Here's the project source tree



References

Resources (data/pretrained model)

- POSO challenge https://challengedata.ens.fr/fr/challenge/33/predisez_la_reponse_attendue.html
- ANSM repository <http://agence-prd.ansm.sante.fr/php/ecodex/index.php>
- Fasttext model <https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md>

NN Papers / Blogs

- Convolutional Neural Networks for Sentence Classification (**REF1**)
 - <https://arxiv.org/abs/1408.5882>
- Text classification with sequential CNN github project (**REF2**)
 - <https://richliao.github.io/supervised/classification/2016/11/26/textclassifier-convolutional/>
- Sensitivity Analysis of Convolutional Neural Networks for Sentence Classification (**REF3**)
 - <https://arxiv.org/abs/1510.03820>
- Text classification blog with LSTM (**REF4**)
 - <https://www.kaggle.com/kredy10/simple-lstm-for-text-classification>

Tools

- Spell checker <https://github.com/barrust/pyspellchecker/blob/master/docs/source/quickstart.rst>
- Scikit-learn <http://scikit-learn.org/stable/>
- XGBoost <https://xgboost.readthedocs.io/en/latest/>
- Keras <https://keras.io/>
- Tensorflow <https://www.tensorflow.org/>
- NLTK <https://www.nltk.org/>