

# **BIAS EVALUATION FOR MODELS OF TOXIC CONVERSATIONS CLASSIFICATION**

---

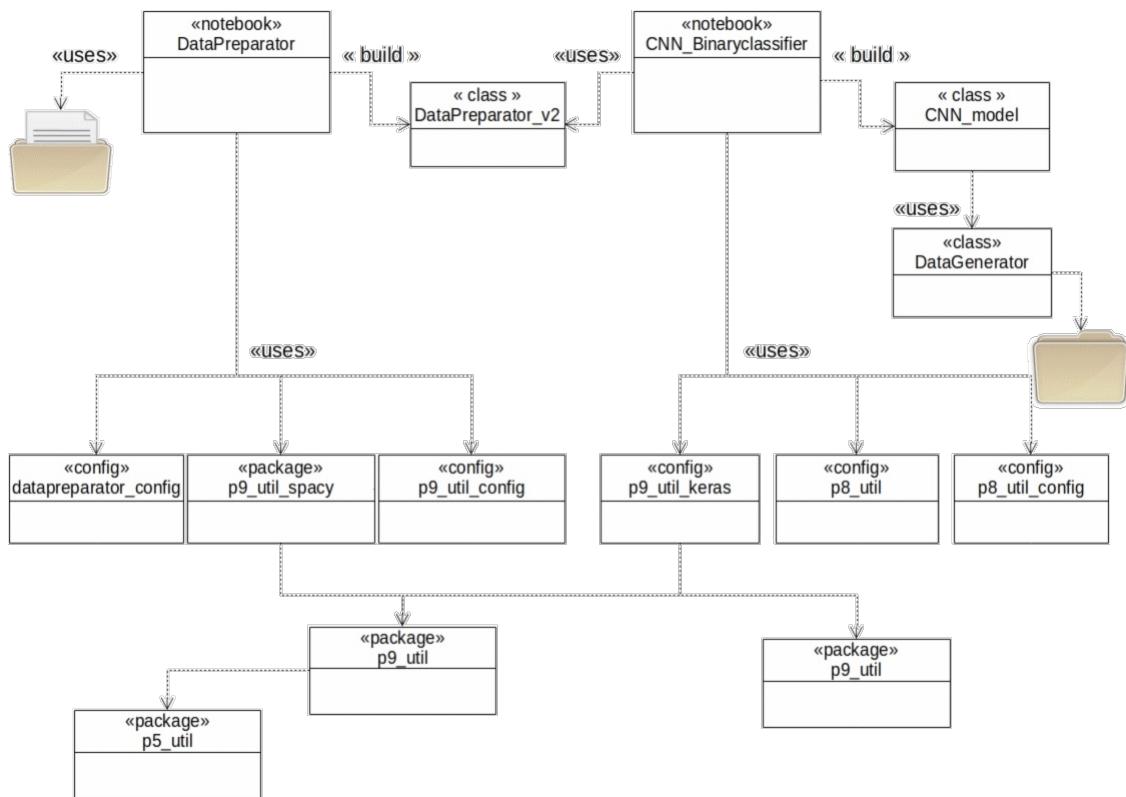
## **Abstract**

This project proposes to implement, as part of a Kaggle competition, a metric to evaluate a machine learning algorithm operating on a population composed of individuals classified in subgroups. These subgroups mark a social or ethnic identity of the individuals making up the population.

Details of projects : <https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification>

- An exploratory analysis of the corpus of texts has been conducted allowing to argue over the choice of learning model.
- Data-preparation process is described and leads to a digital representation of the input data to feed machine learning algorithm.
- The study of machine learning hyper-parameters is conducted and parameters used for building learning model are selected.
- Selected model is trained and results in term of binary classifications are exposed along with bias computation in compliance with formula described in <https://www.kaggle.com/c/>

# Software architecture



This software architecture is drawn with UML syntax. It is organized as layers of packages issued from other Data Sciences projects. They are reused here to face effectiveness, productitiv and quality issues.

All packages are implemented with Python 3.6 language.  
Architecture semantic is interpreted as following :

- The notebook **DataPreparator.ipynb** uses file where dataset is stored in to read it.
- The notebook **DataPreparator.ipynb** builds a data model, named **DataPreparator\_v2.dll**. This an object issued from Python class **DataPreparator\_v2.py** . It contains all configurations and functions (attributes and methods) used in order to produce a dataset ready

for feeding a machine learning algorithm.

- The notebook **CNN\_BinaryClassifier.ipynb** uses **DataPreparator\_v2** Python object in order to build a Keras CNN model for binary classification. This model is saved on harddisk into a **H5** format.
- The **CNN model** uses an instance of **DataGenerator.py** class in order to access data in a bulk way. Data is stored in a folder of hardisk. This allows to save memory ressources, using hardisk as a memory ressource extension and also to train model in a K-fold manner.

# Exploratory analysis

## Features description

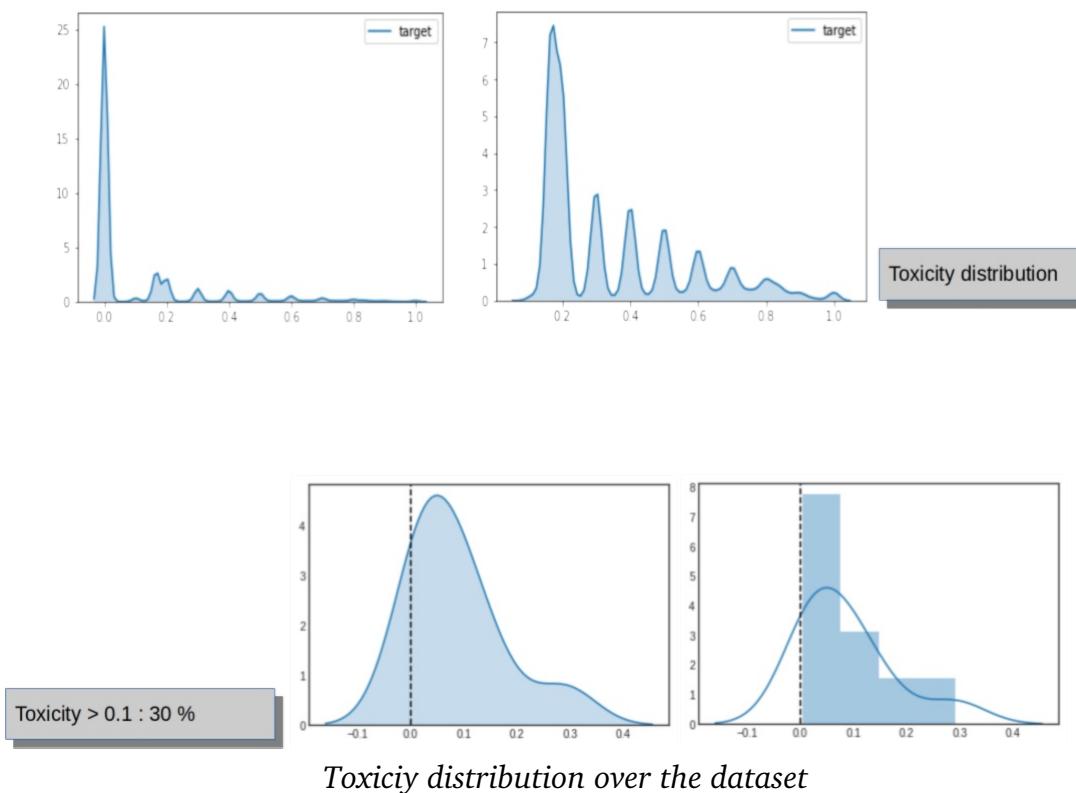
Content & target	
target	0
comment_text	
Toxicity indicators	
severe_toxicity	0
obscene	0
identity_attack	0
insult	0
threat	0
funny	0
wow	0
sad	0
likes	0
disagree	0
sexual_explicit	0
Misc informations	
id	0
created_date	0
publication_id	0
article_id	0
rating	0
identity_annotation_count	0
toxicity_annotation_count	0
parent_id	778646
Identities	
asian	1399744
atheist	1399744
bisexual	1399744
black	1399744
buddhist	1399744
christian	1399744
female	1399744
heterosexual	1399744
hindu	1399744
homosexual_gay_or_lesbian	1399744
intellectual_or_learning_disability	1399744
jewish	1399744
latino	1399744
male	1399744
muslim	1399744
other_disability	1399744
other_gender	1399744
other_race_or_ethnicity	1399744
other_religion	1399744
other_sexual_orientation	1399744
physical_disability	1399744
psychiatric_or_mental_illness	1399744
transgender	1399744
white	1399744

*Dataset features identification*

The dataset is composed of 45 columns described as following:

- Toxicity indicators : these features indicates the type of toxicity
- Identities : these features indicates presence into comments of words or expressions related to an identity. Bias of predictive model will be evaluate against some of these identities.
- Misc informations : these features provide additional informations over comments.
- Content & target : Content are used to feed machine learning model (after a data-preparation and a digitalization process) and target are used as labels to train models.

## Target distribution



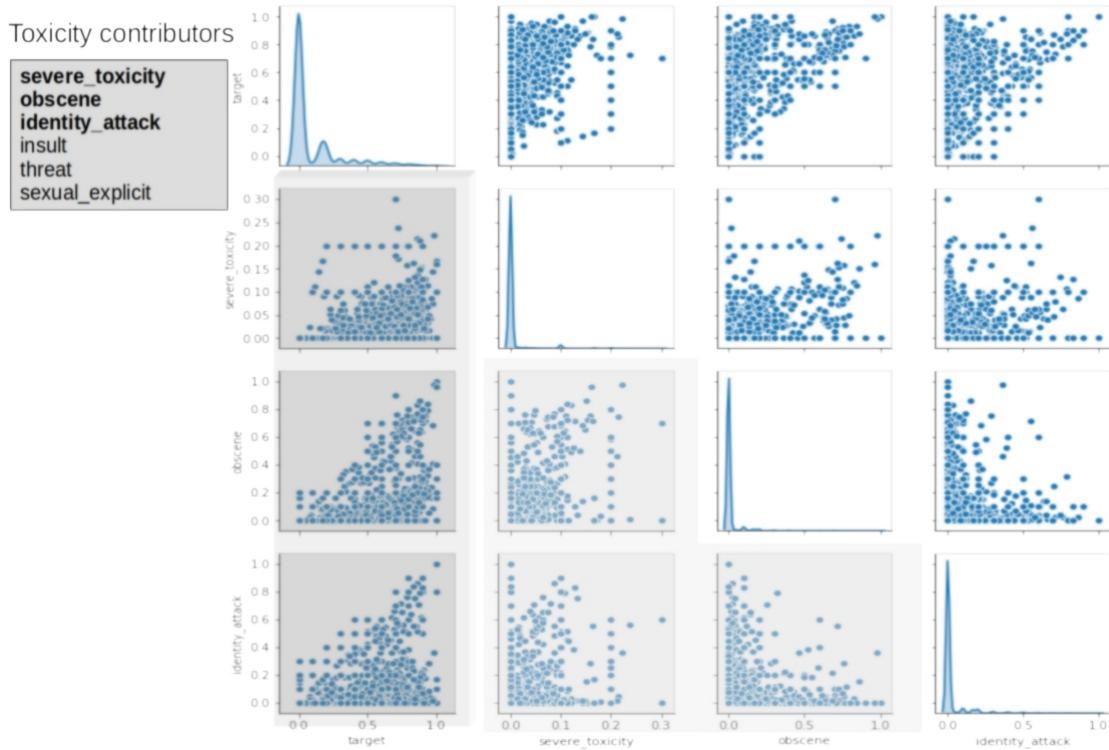
Toxicity values (named target) are ranged from 0.0 (safe comment) to 1.0 (more toxic comment).

Toxicity distribution shows unbalanced classes. Around 30% of comments have a toxicity value greater than 0.0.

Note also that continuous values of toxicity are difficult to interpret and perceive distinctly for a human. E.g it is not clear to perceive

how far a comment with toxicity value of 0.38 is fare from another comment with toxicity value fixed to 0.35.

## Toxicity contributors analysis

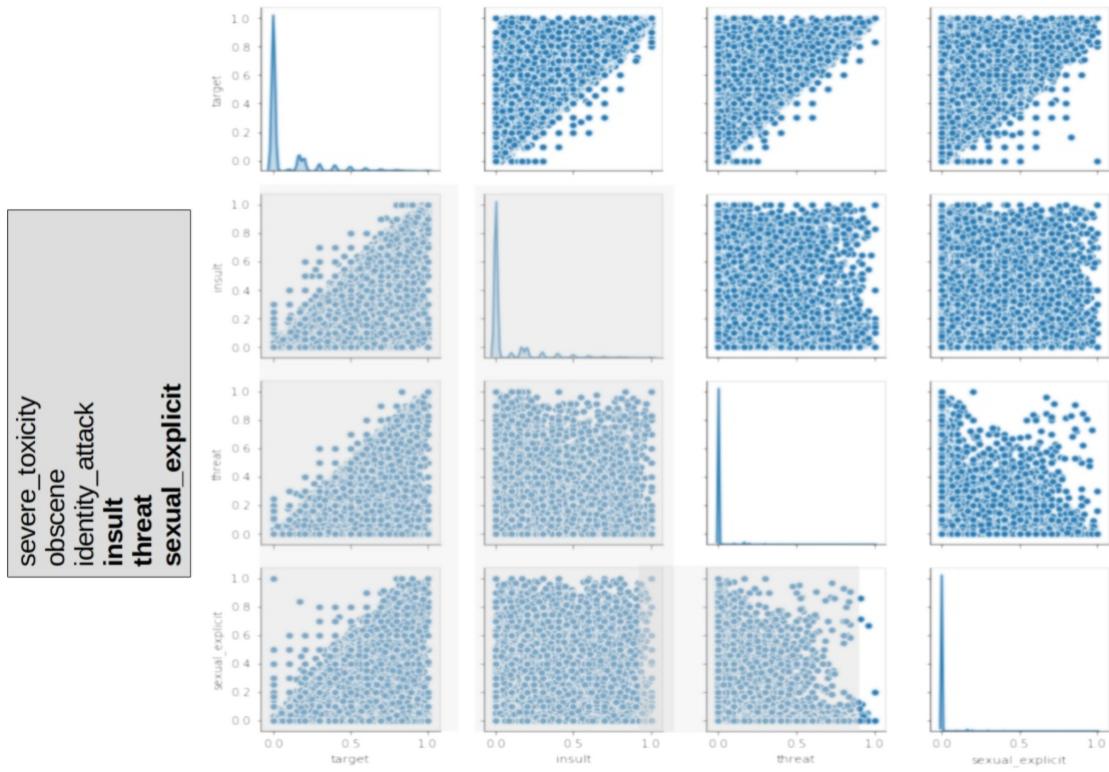


*Multi-variate analysis for toxicity contributors (1)*

The group of features, named **Toxicity distribution** are analysed against target values in order to highlight associations between dependant variables and features.

This multivariate analysis shows that **identity\_attack**, **obscene** and **severe\_toxicity** in extruded area show a linear border of the shape of distributed points along with target. The border of the shape *increases linearly* with toxicity range of values. The fact that a dense area of points appears under the diagonal border let think that contrubution to toxicity may be due to a combinaison ot this three features.

We note also that in non extruded area, multivariate analysis between features **identity\_attack** and **obscene** expose the same kind of linearly border shape that *linearly decreases*



*Multi-variate analysis for toxicity contributors (2)*

On the display above, the multivariate analysis between features **insult**, **threat**, **sexual\_explicit** and **target**, over the extruded area, shows a more accentuated and dense distribution of points below the diagonal border that *linearly increases* with target values.

## ANOVA over identities

QUESTION: is it possible to infer a relation between toxicity and identities ?

Identities are used in order to evaluate bias. Idea is to evaluate the contribution of identities over the toxicity variance

For doing so, ANOVA is proceeded over each one of the groups formed with identities.

Identities values are continuous values ranged from 0.0 to 1.0. These values are reworked and are labelized with values 0,1,... 10. Such transformation allows to consider identities as formed with qualitative groups. These labelized ranges of values may be interpreted as levels

of the considered identity.

The labeling scheme is as following :

Range of group values	Label
0	0
]0.0 , 0.1]	1
]0.1 , 0.2]	2
]0.2 , 0.3]	3
]0.3 , 0.4]	4
]0.4 , 0.5]	5
]0.5 , 0.6]	6
]0.6 , 0.7]	7
]0.7 , 0.8]	8
]0.8 , 0.9]	9
]0.9 , 1.0]	10

*Labeling scheme for identities leading to identity levels.*

E.g. for female identity, such transformation leads to consider 11 groups, ranged from 0 to 10, with statistics values describe below for female identity :

	N	Mean	SD	SE	95% Conf. Interval
<b>Female</b>					
<b>level_0</b>	2782	0.132164	0.216442	0.004104	0.124120 0.140209
<b>level_6</b>	2782	0.160483	0.225222	0.004270	0.152112 0.168854
<b>level_8</b>	2782	0.166551	0.218262	0.004138	0.158438 0.174663
<b>level_7</b>	2782	0.166963	0.220442	0.004179	0.158770 0.175156
<b>level_2</b>	2782	0.167104	0.239145	0.004534	0.158216 0.175993
<b>level_10</b>	2782	0.169066	0.223008	0.004228	0.160778 0.177355
<b>level_5</b>	2782	0.169404	0.230365	0.004368	0.160842 0.177966
<b>level_3</b>	2782	0.170108	0.236792	0.004489	0.161308 0.178909
<b>level_9</b>	2782	0.178096	0.227038	0.004304	0.169658 0.186535
<b>level_4</b>	2782	0.179424	0.243967	0.004625	0.170356 0.188491
<b>level_1</b>	2782	0.186221	0.232819	0.004414	0.177568 0.194874

*Variability of toxicity means along with level of female identity.*

Means values are the means of toxicity computed for each group inside comments marked with female identity.

Question : are differences between these means significant ? Is there an inferable relation between identities labels, ranged from 0 to 10, and toxicity level of a comment ?

For any of the groups in the list above, Levene test is conducted. This test allows to validate homodestaticity hypothesis : variance in each group has equal value. This is the **H0 hypothesis stating over variance equality** between the groups of an identity. If Levene test is significant, considering p-value, then null hypothesis will be accepted and homodestaticity hypothesis will be valid.

Results are produced on the figure below:

Identities	Sampling per group	Levene		ANOVA		
		F-stat	p-value	R <sup>2</sup>	F-stat	p-value
male	3473	23	0.0	0.00	18.53	0.00
female	2782	10.5	0.0	0.00	10.20	0.00
homosexual_gay_or_lesbian	708	13	0.0	0.05	42.54	0.00
christian	2218	36	0.0	0.01	36.02	0.00
jewish	<b>214</b>	<b>1.48</b>	<b>0.14</b>	<b>0.01</b>	<b>2.805</b>	<b>2.e-3</b>
muslim	1045	19.4	0.0	0.03	33.82	0.00
black	780	13	0.0	0.06	52.86	0.00
white	835	14.5	0.0	0.06	56.35	0.00
psychiatric_or_mental_illness	576	14.5	0.0	0.02	10.76	0.00

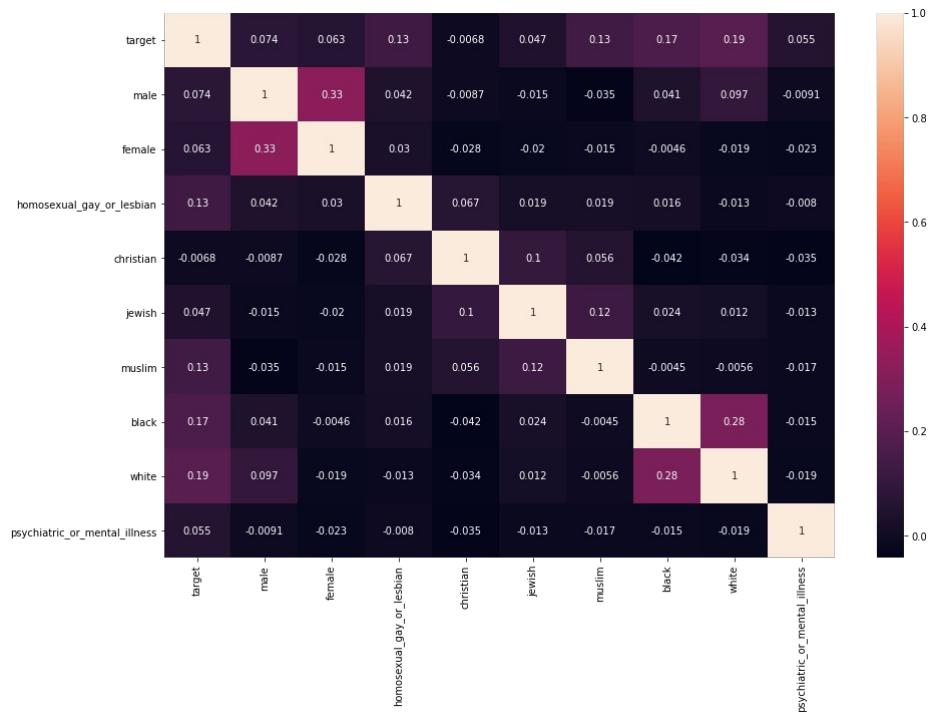
*Analysis of variance with Levene hypothesis tests.*

Levene test seems to be not significant for all identities except jewish

group. Nevertheless, for making F-stat robust considering **normality** and **homoscedasticity hypothesis**, all groups inside each one of the identities have been sampled with the same number of observations. Then, for each identity, considering F-stat magnitude (that measures a ratio of the variability of the mean inside groups and the variability of the mean between groups) and associated p-value << 5.e-2 (that indicates how significant is the result of F-stat), test over differences of means between groups are **regarded as significant**.

### As a conclusion of ANOVA :

- Considering F-stat values for all identities, there is a significant difference between means of groups, groups organized as levels of an identity.
- Considering R<sup>2</sup> values for all identities, explained variance is very weakly for any identity, leading to the conclusion that none of the identities have a size effect over comments toxicity. Correlation matrix here-under partially illustrates this claim. This correlation matrix has been built over the whole observations. We note a symmetric aspect of correlations between couples of features (male, female) and (black, white). These features show a **middle size effect** for explaining target variance.



*Correlation matrix of identities levels*

In addition, diagram below exposes the distribution of a sample of 5000 points for identities that are taken into account for bias evaluation along with toxicity. There is no intuitive evidence of informational shape as pointed for toxicity contributors features on previous section.



*Multivariate analysis of identities levels against toxicity magnetude*

For jewish identity, having the lowest number of observations among all other identities, the F-statistic and related p-value, respectively 2.8 and 2.e-3 shows that there is a significant difference between means of groups inside jewish identity.

Table below shows, for jewish identity, means and standard deviation of any of the groups built on labeled indices ranges of values. For the need of test reliability, all groups have been sampled with the same number of observations, 214.

Jewish identity						
	N	Mean	SD	SE	95%	Conf. Interval
level_0	214	0.137278	0.225208	0.015395	0.107034	0.167523
level_1	214	0.146550	0.204520	0.013981	0.119083	0.174016
level_5	214	0.165121	0.207500	0.014184	0.137254	0.192987
level_2	214	0.167728	0.213974	0.014627	0.138992	0.196464
level_3	214	0.168272	0.203585	0.013917	0.140931	0.195613
level_4	214	0.168771	0.221208	0.015121	0.139063	0.198479
level_7	214	0.177465	0.217144	0.014844	0.148303	0.206626
level_6	214	0.180769	0.199595	0.013644	0.153964	0.207574
level_10	214	0.199158	0.224799	0.015367	0.168968	0.229348
level_8	214	0.209509	0.229285	0.015674	0.178716	0.240301
level_9	214	0.218742	0.237592	0.016241	0.186834	0.250649

*Means of toxicity along with identities groups formed with identities levels*

Means show differences in between labelized groups. Are these differences significant? If yes, then it will be allowed to conclude that such groups, based on jewish identity, do explain variancy of toxicity. In that case, it will be stand that a relation between toxiciy level and jewish identity exists.

Table below shows the result of ANOVA over jewish identity.

	Coef	std err	t	P> t	[0.025	0.975]
Intercept	<b>0.1373</b>	<b>0.015</b>	<b>9.251</b>	<b>0.000</b>	<b>0.108</b>	<b>0.166</b>
C(jewish)[T.level_1]	0.0093	0.021	0.442	0.659	-0.032	0.050
<b>C(jewish)[T.level_10]</b>	<b>0.0619</b>	<b>0.021</b>	<b>2.949</b>	<b>0.003</b>	<b>0.021</b>	<b>0.103</b>
C(jewish)[T.level_2]	0.0304	0.021	1.451	0.147	-0.011	0.072
C(jewish)[T.level_3]	0.0310	0.021	1.477	0.140	-0.010	0.072
C(jewish)[T.level_4]	0.0315	0.021	1.501	0.134	-0.010	0.073
C(jewish)[T.level_5]	0.0278	0.021	1.327	0.185	-0.013	0.069
<b>C(jewish)[T.level_6]</b>	<b>0.0435</b>	<b>0.021</b>	<b>2.072</b>	<b>0.038</b>	<b>0.002</b>	<b>0.085</b>
<b>C(jewish)[T.level_7]</b>	<b>0.0402</b>	<b>0.021</b>	<b>1.915</b>	<b>0.056</b>	<b>-0.001</b>	<b>0.081</b>
C(jewish)[T.level_8]	0.0722	0.021	3.442	0.001	0.031	0.113
C(jewish)[T.level_9]	0.0815	0.021	3.882	0.000	0.040	0.123

*Analysis of variance of means of toxicity along with identities groups formed with identities levels*

ANOVA has been conducted over 214 observations and shows :

- $R^2$  value, 1.e-2 is a weak value. Correlation between toxiciy and identity is weak. There is no **linear relation** between jewish identity and toxicity. Correlation matrix here-under supports this conclusion :
- From the column  $P>|t|$  the groups labeled 6, 7, 10 have significant difference with intercept, that is group level 0. Post-hoc test will allows to compare groups level 6, 7 and 10 in between each-other.

## Conclusions over features analysis

Features have been splitted into two categories : contributors to toxicity of comments and identities. It has been shown that contributors to toxicity has a strong effect over the toxicity magnitude, whereas identities have a weak effect over toxicity. Unintended bias in machine learning classification comes with the fact that despite that last point, comments tend to be classified as toxic when they embedd some identities terms or expressions.

Absence of identity effects over toxicity also let think that the process involved in toxicity and identity metrics can be safely taken into account. This lead to consider as reliable comments labelization used in a supervised machine learning model.

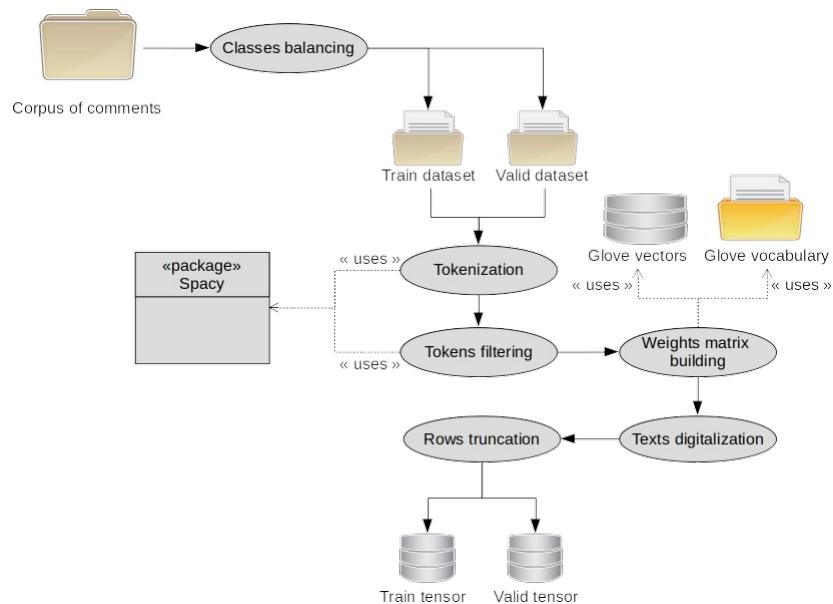
Same remark also stands for toxicity contributors. Their co-variancy behaviour along with toxicity is compliant with the intuitive sens, stating that a comment that embeds insulting terms or expressions is welling to render it toxic.

Proper algorithms should be able to take into account the presence of toxicity contributors along with identities in comments in order to proceed to classification.

Also, post-hoc tests shows no evidence of an existance of treshold effect identifying differences between means of different groups inside an identity.

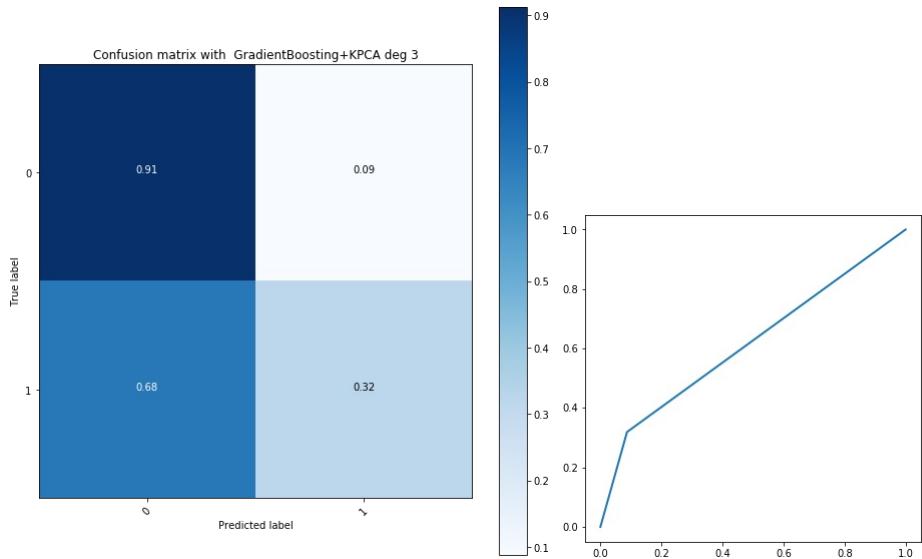
# Data preparation & digitalization process

This process allows to transforms a corpus of comments into a digitalized representation, allowing to feed an algorithm. The global process is described below :



*Software engineering involved into data-preparation and digitalization process*

- Spacy package is used for NLP processing. Using it, sequences of words are preserved using.
- Due to toxicity distribution study in previous sections, balancing dataset is required in order to ensure toxic and non toxic comments to be processed "equally". Otherwise, contributions to cost function will be mainly due to non toxic comments that will lead to a biased prediction model. The picture above do represents a model obtained with Gradient boosting algorithm without dataset balancing. All non toxic comments are very well classified while predictions for toxic comments are not better then random prediction leading to weak performances as shown over the RAUC curve.



*Gradient boosting algorithm binary classification performances with unbalanced classes*

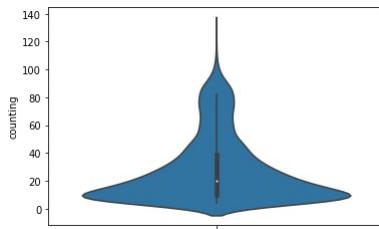
- As the result of the digitalization process, a tensor with 3 dimensions is produced, for train and validation dataset. It is represented as the figure below :

		Dim 0	...	DIM D
Text t	Token 0	$\text{coeff}_{t,0,0}$		$\text{coeff}_{t,0,D}$
	Token 1	$\text{coeff}_{t,1,0}$		$\text{coeff}_{t,1,D}$
	...			
	Token N	$\text{coeff}_{t,N,0}$		$\text{coeff}_{t,N,D}$
Text $t+1$	Token 0	$\text{coeff}_{t+1,0,0}$		$\text{coeff}_{t+1,0,D}$
	Token 1	$\text{coeff}_{t+1,1,0}$		$\text{coeff}_{t+1,1,D}$
	...			
	Token N	$\text{coeff}_{t+1,N,0}$		$\text{coeff}_{t+1,N,D}$

Dimension 1 : Number of texts    Dimension 2 : Max text length    Dimension 3 : embeddings

*3D tensor representation of digitalized dataset*

- Texts truncation : all tokenized texts should have the same number of tokens. This allows to feed NN algorithm with **same embedding layer size**. This operation takes place after the filtering process, for keeping the model with usefull informations. The max length size for the number of tokens has been fixed to 100. This value comes with the text word length distribution described on diagram below :



*Text length distribution over the dataset*

Trunc of texts after digitalization allows to drive this process based on an objective criteria, such as magnitude of the vectors. Vectors with the smallest magnetude will be pushed out of digitalization representation. While doing so, relevant information that will contribute to decrease of loss function will be kept.

- \* Texts padding : digitalized texts with number of tokens less then max length will be padded with zero vector.

## Word embeddings

Glove allows to use different kind of model language for Naturel Language Processing. The one used here is `en_core_web_lg`, in which, each word in vocabulary is represented as a vector in a 300 dimensions space. Words vectors have been built using web texts and comments issued from variou social networks.

A dictionary structured as `{word:glove_coefficient}` is built from Glove source.

Once built, dictionary allows to build a vector for every word in vocabulary issued from tokenizer.

Endly, weights matrix is built from vocabulary issued from tokenizer.

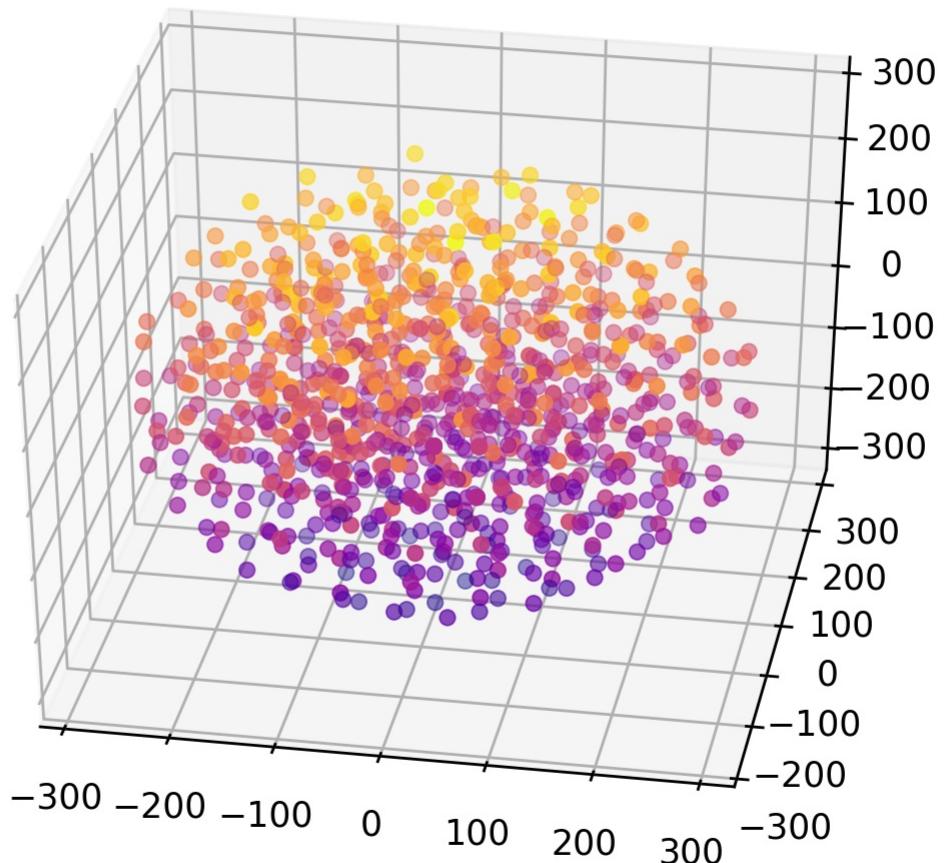
Such process is summarized with sequences here-under :

- `dict_glove_word_coeff <- processing Glove file name`
- `vocabulary_word, index <- tokenizer`
- `weight_vector = dict_glove_word_coeff[vocabulary_word]`
- `weight_matrix[index] = weight_vector`

# Features correlations

Statistic analysis as shown that some features expose linear correlations with toxicity values.

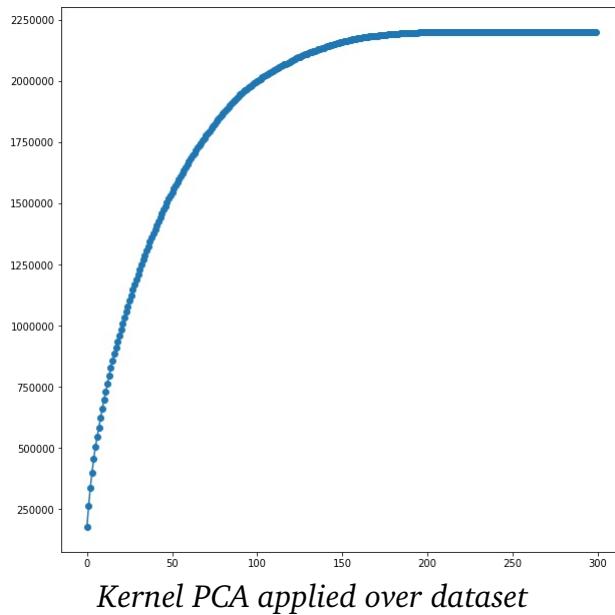
Figure below shows 1000 digitalized comments with process described in previous section. t-SNE operator is applied over 3 dimensions. There is no evidence that an hyperplan exists that is able separate groups of points.

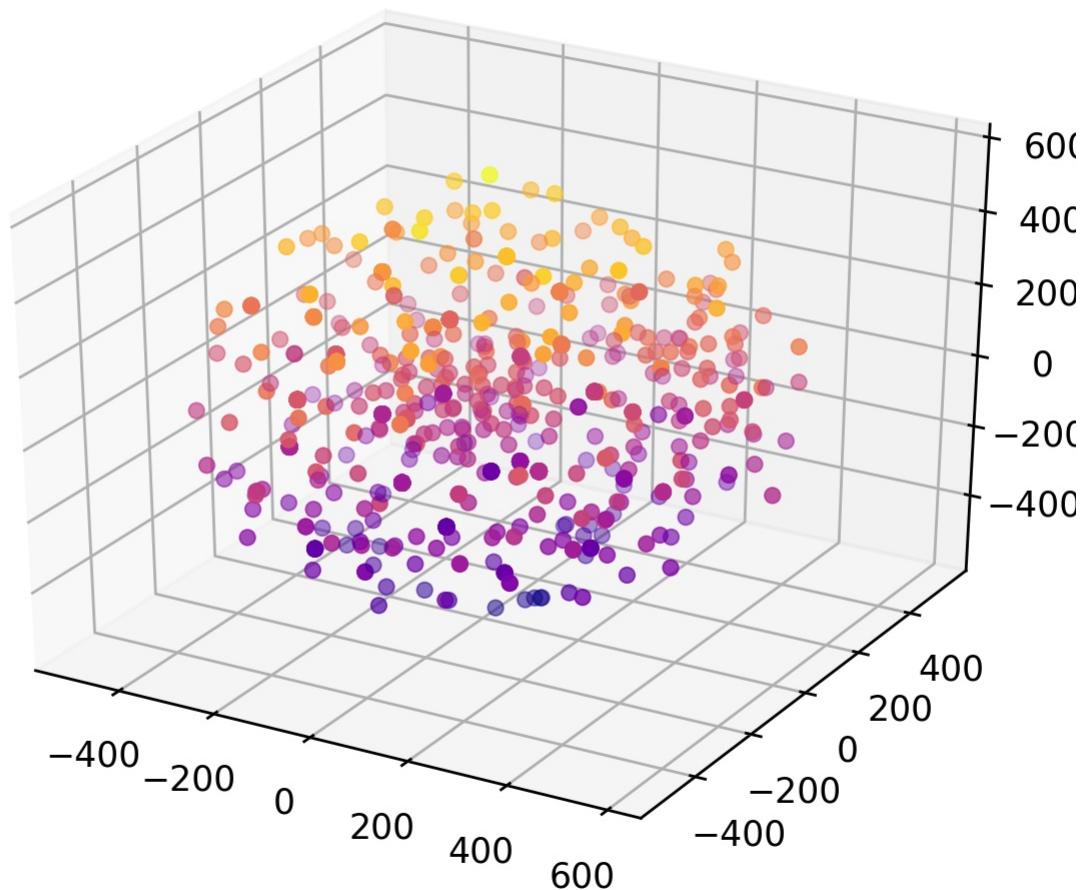


*3D t-SNE reduction over dataset with PCA operator*

On figure below, a kernel PCA operator has been applied with kernel

of order 5. The kernel trick allows to process a non linear model as a linear model in an Hilbert space of larger size then the original space.





*3D t-SNE reduction over dataset with order 5 Kernel PCA operator*

## CNN classifier

From statistic analysis conducted in previous sections, it is shown that the presence of some type of words or expressions in a comment is correlated to the toxicity level. The problem can be formulated as the way to identify in a comment, what has been named toxicity contributors and identities and relate them to a level of toxicity. Due to that, the expected prediction model should be able to learn such relations connecting identified structures with toxicity level.

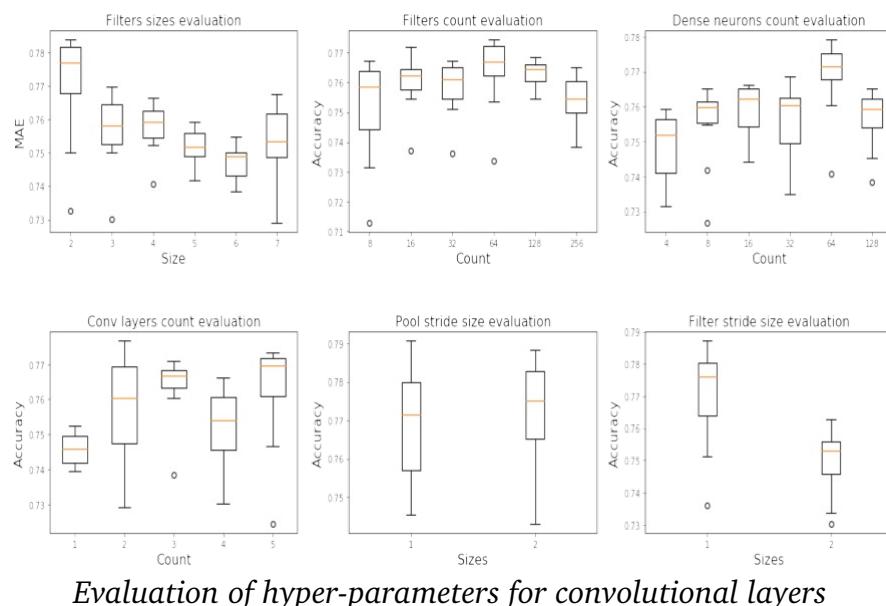
The ground problem is focused on textual structures identification.

CNN architectures are famous and relevant to accomplish a such task. In addition, a bias metric defined in <https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/overview/evaluation> will be applied over the model in order to evaluate how fare the model is able to distinguish toxic comments from non toxic comments that embeds identities terms or expressions.

Rather than predict toxicity level thanks to a regression algorithm, a classifier is used to complie with the model bias evaluation.

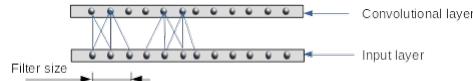
## CNN hyper-parameters selection

### CNN hyper-parameters selection for convolutional layers



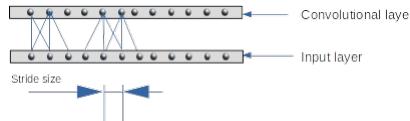
These parameters are related to CNN structure :

- The size of the convolution filters. When processing text, this parameter amounts to setting n-grams for the neighborhood of words. Such a structure in word processing is related to the semantics of words, indicating that words of the same neighborhood have the same meanings.



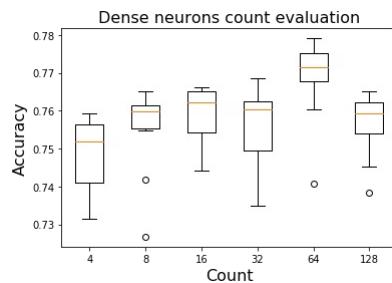
*Illustration of filter size effect over the input layer*

- The number of convolutional filters, leading to features maps. This will give the model the potential for identifying all kinds of textual structures in the corpus of comments. A single feature map identifies the same structures that may be repeated in a comment.
- The number of convolutional layers.
- The filter stride size. This parameter fixes the window stride for N-GRAM structures detection.



*Stride size effect over the input layer*

## CNN hyper-parameters selection for dense layers



*Evaluation of the number of neurons into dense layer*

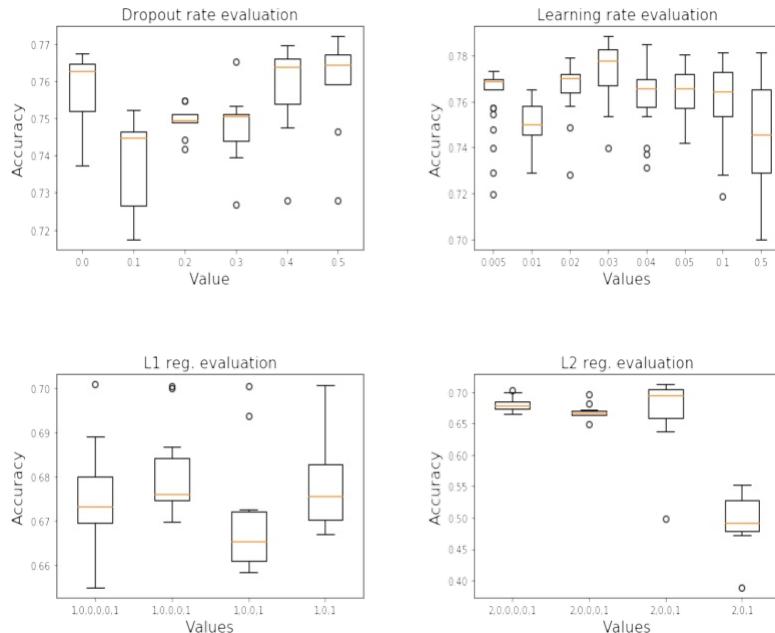
These parameters are related to dense layers structure and they are:

- The number of dense layers. It has been fixed to 1.
- The number of neurons in the dense layer. Dense layers after convolutional layers allows to classify features issued from convolutional layers.

## Algorithms selection :

- Gradient descent algorithm is RMSprop.
- Loss function algorithm is binary cross entropy.

## CNN hyper-parameters selection for optimization :



*Hyper-parameters evaluation for model optimization*

- Dropout rate. This parameter controls the model complexity and hence, its capacity to generalize.
- Learning rate. This parameter controls the speed of the learning.
- L1 regularization. This parameter controls the complexity of the algorithm, hence its ability to generalize. This regularization process is applied to the cost function. Complexity regularization is achieved while killing some features in the learning process. L1 regularization is known to be less stable than L2 regularization, because loss function convexity is impacted with L1 terms.
- L2 regularization. This parameter controls the complexity of the algorithm, hence its ability to generalize. This regularization process is applied to the cost function.

# Classification results

## Binary classification

This step has been performed by splitting toxicity values into 2 classes :

- $0.0 \leq \text{toxicity} < 0.5$  : safe comments
- $0.5 \leq \text{toxicity} \leq 1.0$  : toxic comments

## Data generators

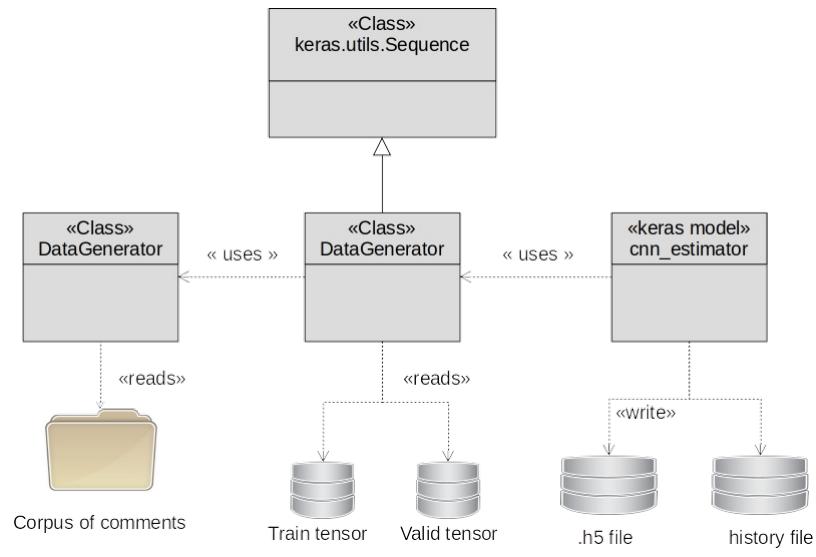
Because of memory resources issues, DataGenerator objects have been created and used for both, training and validation operations.

Such objects allow to pump bulk of data stored on harddisk.

Architecture below shows how the model works :

DataGenerator class inherits from keras.utils.Sequence. This allows Keras model to use object issued from this class as a source of data, pumping it iteration after iteration in the train or validation operations, until all data partitioned into files over hard disk have been processed.

DataGenerator is created using DataPreparator object. Such objects contain all operations for data preparation and digitalization process along with configuration parameters used inside these process. This architecture allows an **automation digitalization process** to take place.

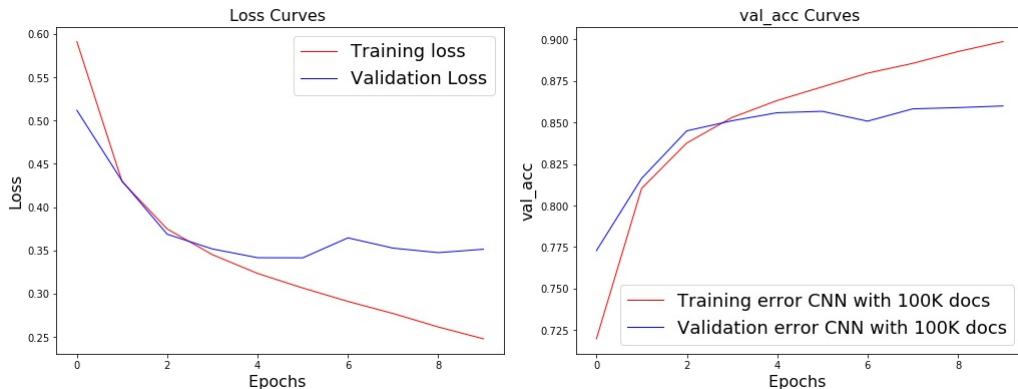


*Software architecture for the learning process*

Depending on object configuration, data sources may be train dataset or validation dataset. A callback function, issued from object of type `keras.callbacks.ModelCheckpoint` allows to save the best model issued from train operation.

## Training / validation performances

The number of trainable parameters is closed to 715K Model tends to overfit, despitely dropout rate of 0.3



*Accuracy and loss functions performances for submission model*

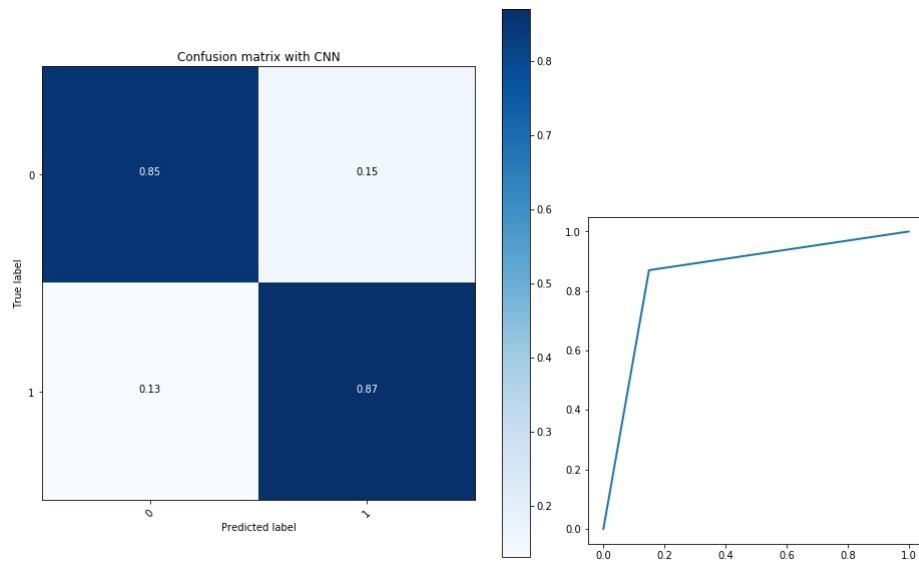
## Classification performances

Results below have been obtained with training near to 100K comments and validated over closed to 10K comments.  
Confusion matrix below shows that model is able to properly classify

:

- 85% of non toxic comments
- 87% of toxic comments

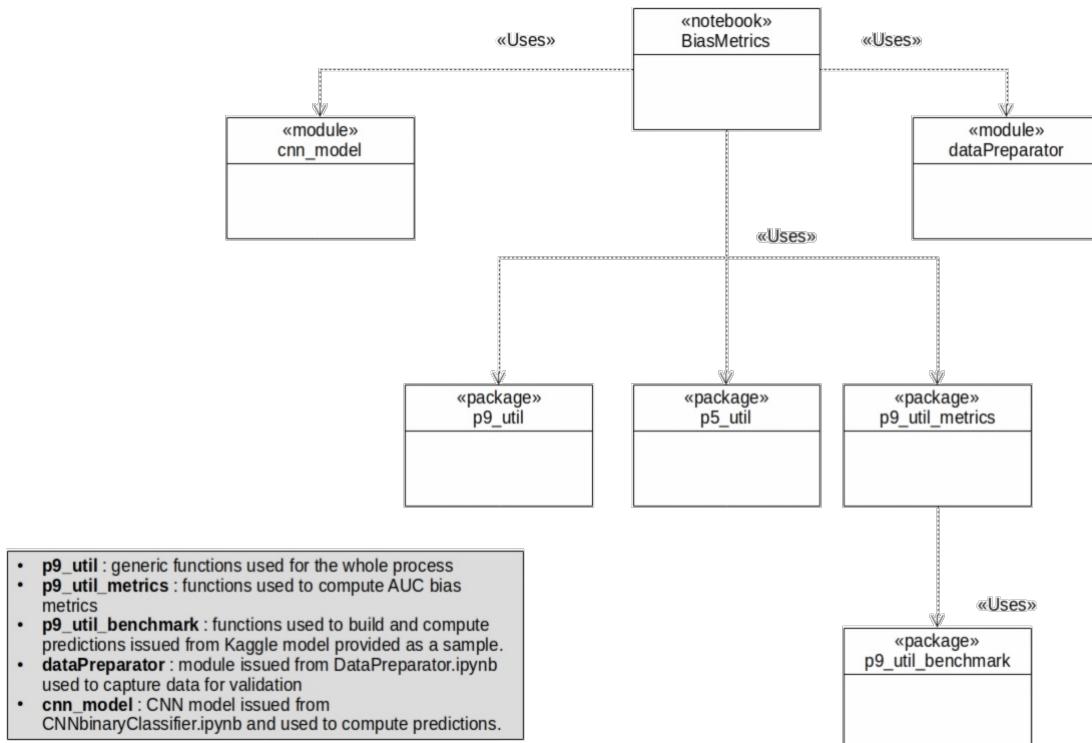
Area Under RO Curve is AUC=0.86



*Confusion matrix and ROC curve for the submission model*

## Bias evaluation

Graphic below shows software architecture involves in the bias evaluation process.



# Bias evaluation of submission model

Submission model is the one that aims to be submitted to the Kaggle competition. This is the one that has been described and built in this report.

A bias metric is defined in order to evaluate AUC over predictions, thus, for each one of the identities. For each identity, False Positive (FP) rate and True Positive (TP) rate are jointly evaluated by defining area under ROC curve.

This metric measures three aspects of the classification for each identity :

- **BNSP\_AUC** : the ability of the model to distinguish safe comments without any reference (words, expressions) to the identity (such comments belongs to the group named **Background Negative**) from the toxic comments with reference to the identity (such comments belongs to the group named **Subgroup Positive**).
- **BPSN\_AUC** : the ability of the model to distinguish toxic comments without any reference (words, expressions) to the identity (such

comments belongs to the group named **Background Positive**) from the safe comments with references to the identity (such comments belongs to the group named **Subgroup Negative**).

- **SUBGROUP\_AUC** : the ability of the model to distinguish toxic comments from non toxic comments, all comments having references to identity.

While doing so, the classification model performance is scanned for each identity. Table below shows these three values, identity per identity.

	<b>bnsn_auc</b>	<b>bpsn_auc</b>	<b>subgroup</b>	<b>subgroup_auc</b>	<b>subgroup_size</b>
<b>2</b>	0.837397	0.669649	homosexual_gay_or_lesbian	0.645546	149
<b>6</b>	0.880691	0.641463	black	0.661654	199
<b>7</b>	0.888992	0.652551	white	0.680226	335
<b>5</b>	0.867801	0.713312	muslim	0.719677	222
<b>8</b>	0.917395	0.727446	psychiatric_or_mental_illness	0.784848	43
<b>0</b>	0.852952	0.810845	male	0.802190	356
<b>1</b>	0.853437	0.821859	female	0.813832	404
<b>4</b>	0.838606	0.852951	jewish	0.831250	68
<b>3</b>	0.841626	0.851564	christian	0.832449	269

*Identities AUC for the submission model*

Considering measures of **SUBGROUP\_AUC**, the model confuses to distinguish toxic comments from safe one for those comments with **homosexual\_gay\_or\_lesbian** identity. For this identity, probability that model proceeds to a proper distinction between toxic and non toxic comments in less than 70%.

This probability of distinction is nearly the same value, 70%, for toxic comments without any reference to homosexual, gay or lebian and safe comments that embedds these references.

The final score is computed thanks to the overall AUC that is 0.86 and the three aspects of AUC computed separately for each identity, with the formula :

$$score = 0.25 * AUC_{global} + 0.75 * \left( \sum_{g=1}^{N_g} \frac{1}{N_g} (bsnphauc_g + bsnphauc_g) + subgroupa \right)$$

Where :

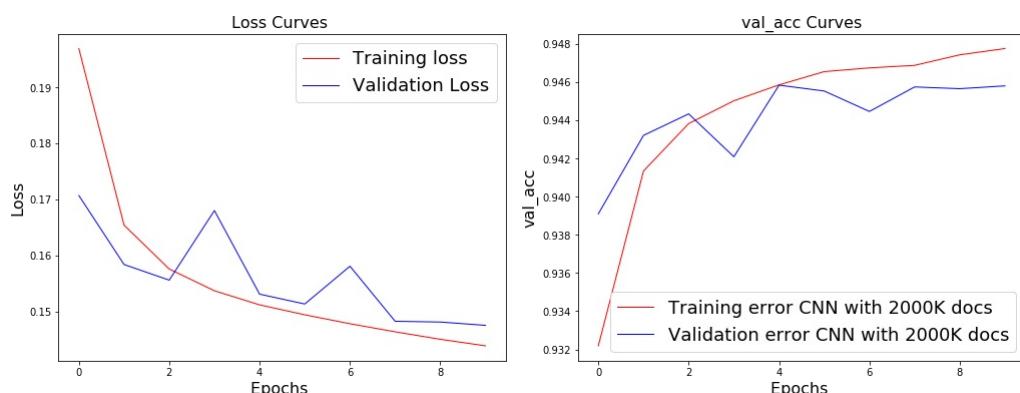
- p is the power, fixed to 5
- Ng is the number of identities involved in the bias evaluation, fixed to 9
- AUCglobal is the AUC value computed in previous section, with value = 0.86
- Other AUC values are the one extracted from bias table of identities.
- 0.25 and 0.75 are weights factors that favor ability of the model to distinguish the toxicity of texts separately for each identity over global separability ability.

For the model to be submitted to Kaggle, this formula comes with value 0.79 This may be interpreted as a bias value of 20%

## Bias evaluation of benchmark model

Benchmark model is the one provided on Kaggle web site, in order to expose implementation of bias algorithm.

This is a CNN model trained with close to 2 millions comments, say, the whole dataset. It holds close to 160K parameters.



Accuracy and loss functions performances for benchmark model

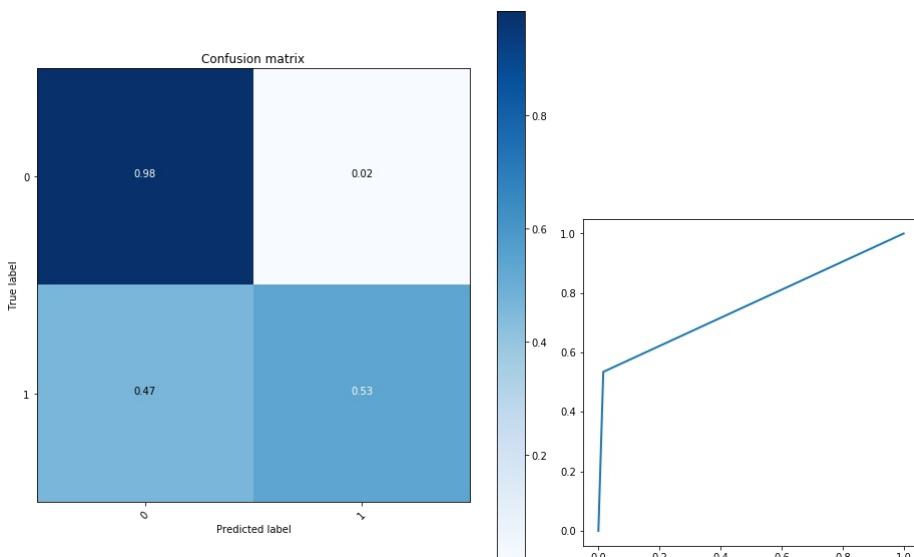
Bias table below shows AUC values for identities :

	bnsn_auc	bpsn_auc	subgroup	subgroup_auc	subgroup_size
6	0.959931	0.768903	black	0.807377	3006
7	0.965931	0.773480	white	0.819055	5077
2	0.963686	0.781776	homosexual_gay_or_lesbian	0.820112	2221
5	0.957757	0.819413	muslim	0.845473	4262
4	0.948680	0.857917	jewish	0.862545	1509
8	0.948939	0.856340	psychiatric_or_mental_illness	0.864638	982
0	0.946800	0.875431	male	0.883611	9050
1	0.944072	0.883673	female	0.889211	10671
3	0.937363	0.909032	christian	0.907213	8081

*Identities AUC for the benchmark model*

Final score for benchmark model, measuring the bias model is 0.89

Confusion matrix below shows that, despite better performances for bias metric, benchmark model confuses more, globally, to distinguish toxic comments from safe one, with AUC = 0.75 against 0.86 for our submission model.



*Confusion matrix and ROC curve for the benchmark model*

# Conclusions

Data preparation process should include emoticons.

Features may be increased with POS tags, NER, organizations tags, all such tags provided with Spacy model of language.