

**Code explanations**  
**Elad Kapuza, Shir Shtinits, Hillel Merran**

**SVM**

Function	Description	Parameter	Return
RemoveHashTag	Remove '#' from an Hash tag. This function is used to decompose it	word - a word (string)	the word without '#'
SplitCamelWordIntoList	Split a camel case word into words list e.g HelloWorld into [Hello, World]	camelWord - a camel case word (string)	List of words from the camel case word
CreateCorpusAndPostListsAndTagList	Create a corpus of termes, a list of posts, of tags, from the data posts	<b>corpus</b> (return) - list of all the words in the data <b>postsList</b> (return) - list of posts, each post is a list of words <b>tagList</b> (return) - list of tags, each tag is in {0: non-racist, 1: racist} <b>stopWordsList</b> - a list of stopwords to remove from the posts	-
PostToVector	Convert a post to a vector for the learning algorithm	<b>corpus</b> - list of all words in the data <b>post</b> - the post to convert. the post is a list of words	the vector which represents the post
StringToVector	Convert a post to a vector for the learning algorithm	<b>corpus</b> - list of all words in the data <b>stopWordsList</b> - list of stopwords to remove from the post <b>str</b> - a post, as a string	the vector which represents the post
CreateVectors	Convert a list of posts to a list of vectors	<b>corpus</b> - list of all words in the data	-

		<b>postsList</b> - list of the posts to convert. each post is a list of words <b>vectorsList</b> (return) - list of vectors. each vector represents a post	
main	<ul style="list-style-type: none"> <li>• create a stemmed corpus with words from the data</li> <li>• create a list of all posts and a list of all tags</li> <li>• create the vectors from the posts</li> <li>• run SVM with Cross-Validation (5-folds)</li> </ul>	-	-

## Optimization

Function	Description	Parameters	Return
sum_of_features	Calculate the sum of the different clique features vectors of a sentence	<b>y</b> - tag of a post. $y = (\text{post\_tag}, [s1\_tag, \dots, sK\_tag])$ <b>s</b> - a post (list of sentences). $s = [s1\_string, \dots, sK\_string]$ <b>n</b> - length of the features vectors	the sum of the vectors
Hamming_loss	Calculate the Hamming loss of a tag. Namely how many coordinates of 2 binary vectors are	<b>w</b> - a first tag. $w = (\text{post\_tag}, [\text{sentence1\_tag}, \dots, \text{sentenceK\_tag}])$ <b>z</b> - a second tag. $z = (\text{post\_tag}, [\text{sentence1\_tag}, \dots, \text{sentenceK\_tag}])$	the Hamming loss
argmin	Resolve the 4th step in MIRA algorithm. That is calculate the weight vector for iteration $i+1$	<b>w</b> - original weights vector <b>y</b> - real tag of a post. $y = (\text{post\_tag}, [\text{sentence1\_tag}, \dots, \text{sentenceK\_tag}])$ <b>s</b> - a post (list of sentences). $s = [s1\_string, \dots, sK\_string]$	weight vector for next iteration

## Naïve Bayes

Function	Description	Parameters	Return
simplify_sentence	Convert a sentence (string) to a list of stemmed words, without stopwords	<b>row</b> - a sentence (string) <b>stopWordsList</b> - a list of stopwords	the list of the words of the simplified sentence
main	create a list of all posts and tags run Naive-Bayes Classifier with Cross-Validation (5-folds)	-	-

## MEMM

Function	Description	Parameters	Return
simplify_sentence	Convert a sentence (string) to a list of stemmed words, without stopwords	<b>row</b> - a sentence (string) <b>stopWordsList</b> - a list of stopwords	the list of the words of the simplified sentence
clique_score	Calculate the score of a tag for a clique	<b>s1, s2</b> - sentences of the clique (string) <b>post_tag</b> - tag of the post {0: non-racist, 1: racist} <b>tag1, tag2</b> - tags of s1, s2 {-1: anti-racist, 0: neutral, 1: racist} <b>stopWordsList</b> - a list of stopwords	the score of the clique
score	Calculate the score of a tag for a post (and sentences)	<b>weights</b> - vector of weights, one for each feature <b>y</b> = (post_tag, [sentence1_tag, ..., sentenceK_tag]) - tag of the post <b>s</b> = [sentence1, ..., sentenceK] - list of sentences (string) of the post	the score of the clique
max_clique_score	Determine the tag of the first sentence in a clique which have the highest clique score	<b>weights</b> - vector of weights, one for each feature <b>clique_no</b> - index of the clique <b>current_2nd_sentence_tag</b> - tag of the 2nd sentence of the clique <b>sentence1</b> - first sentence of the clique (string) <b>sentence2</b> - second sentence of the clique (string) <b>stopWordsList</b> - list of the stopwords <b>viterbi</b> - dictionary for Viterbi algorithm. {key=state: value=score}	<b>backpointer</b> - tag of the first sentence which maximize the score of the clique <b>score</b> - maximum score of the clique

argmax	Determine the tag of the sentences given the tag of the post using Viterbi algorithm	<b>weights</b> - vector of weights, one for each feature <b>d_tag</b> - tag of the post {0: non racist, 1: racist} <b>s</b> - list of the sentences (string) of the post	list of tags. each tag is in {-1: anti-racist, 0: neutral, 1: racist} and correspond to a sentence
classifier	Determine the tag of the post and sentences given the tag of the post using Viterbi algorithm	<b>weights</b> - vector of weights, one for each feature <b>s</b> - list of the sentences (string) of the post	tag of the post and sentences - (post_tag, [s0_tag, ..., sM-1_tag])
main	create a list of all posts and a list of their tag run the MIRA algorithm to learn the vector of weights test the model with Cross-Validation (5-folds)		

## Features

Function	Description	Parameters	Return
StemWord	Convert a word to its stemmed form	word - a word (string)	The stemmed form of the word
CreateRacistStemmedCorpus	Create a corpus of racist stemmed word from csv file	-	Set of racist words
CreateAntiRacistStemmedCorpus	Create a corpus of anti-racist stemmed word from csv file	-	Set of anti-racist words
RemoveHashTag	Remove '#' from an Hash tag. This function is used to decompose it	word - a word (string)	the word without '#'
SplitCamelWordIntoList	Split a camel case word into words list e.g HelloWorld into [Hello, World]	camelWord - a camel case word (string)	List of words from the camel case word
simplify_sentence	Convert a sentence (string) to a list of stemmed words, without stopwords	<b>row</b> - a sentence (string) <b>stopWordsList</b> - a list of stopwords	the list of the words of the simplified sentence
CreateStemmer	Create stemmer object	-	Stemmer object
Similarity	Calculate similarity value between two sentences by counting common words	<b>sentence1</b> - first sentence (list of words) <b>sentence2</b> - second sentence (list of words)	Similarity value
IsSimilarityLow	Check if the similarity level between two sentences is low	<b>sentence1</b> - first sentence (list of words) <b>sentence2</b> - second sentence (list of words)	True if the similarity level is low and false otherwise
IsSimilarityHigh	Check if the similarity level between two sentences is high	<b>sentence1</b> - first sentence (list of words)	True if the similarity level is high and false otherwise

		<b>sentence2</b> - second sentence (list of words)	
ContainsExclamationMark	Check if the sentence contains exclamation mark	sentence - sentence (string)	True if the sentence contains exclamation mark and false othrewise
BothContainExclamationMark	Check if the both two sentences contain exclamation mark	<b>sentence1</b> - first sentence (string) <b>sentence2</b> - second sentence (string)	True if the both two sentences contain exclamation mark and false othrewise
ContainsRacistWord	Check if the sentence contains a racist word	<b>sentence</b> - sentence (list of words) <b>racistWord</b> - racist word (string)	True if the sentence contains the racist word
ContainsAnyRacistWord	Check if the sentence contains any racist word from the racist words corpus	sentence - sentence (list of words)	True if the sentence contains any racist word
ContainsAntiRacistWord	Check if the sentence contains an anti-racist word	<b>sentence</b> - sentence (list of words) <b>antiRacistWord</b> - anti-racist word (string)	True if the sentence contains the anti-racist word
ContainsAnyAntiRacistWord	Check if the sentence contains any anti-racist word from the anti-racist words corpus	sentence - sentence (list of words)	True if the sentence contains any anti-racist word
Polarity	Calculates the polarity of the sentence (sentiment level between -1 to 1)	sentence - sentence (string)	Polarity value
IsPositiveSentimental	Check if the sentence has a positive sentiment	sentence - sentence (string)	True if the sentence has a positive sentiment and false otherwise
IsNegativeSentimental	Check if the sentence has a negative sentiment	sentence - sentence (string)	True if the sentence has a negative sentiment and false otherwise



Subjectivity	Calculates the subjectivity level of the sentence (between 0 to 1)	sentence - sentence (string)	Polarity value
IsSubjective	Check if the sentence is subjective	sentence - sentence (string)	True if the sentence is subjective and false otherwise
IsObjective	Check if the sentence is objective	sentence - sentence (string)	True if the sentence is objective and false otherwise
HasAdjective	Check if the sentence contains adjectives	sentence - sentence (string)	True if the sentence contains adjectives and false otherwise
clique_to_features	Creates list of features with values of 0 or 1	<b>post_tag</b> - The true tag value of the post (int) <b>s1_tag</b> - The true tag value of a sentence in the post (int) <b>s2_tag</b> - The true tag value of the following sentence in the post (int) <b>s1</b> - A sentence in the post <b>s2</b> - The following sentence in the post <b>stopWordsList</b> - Corpus of stop words	List of features with values of 0 or 1