

## About:

This Project is about taking pairs of sentences or any other combination of pairs (X1.... Xn, Y2.... Ym) and learn the relation between the pairs.

The main object is to return a prediction Y for a new input X Using Neural Networks.

It might be related to a lot of subjects in our lives, you can understand the idea by seeing the datasets we chose to work with.

**The System can work and generate pairs from 3 types of datasets:**

1. Only Positive Pairs (Dynamic Generating Negative samples)
2. Mixed Positive and Negative (Organize Positive and Negative samples).
3. Already Organized Pairs with class. (No change needed)

**The System have 3 pre-processing pipelines to use, from soft preprocessing to hard preprocessing.**

**The System 5 Neural Networks Models, some work better some less, it more depends on the data itself.**

**You can save all the preprocessing files and load to train on any Neural Model you want.**

**The Neural Model also saved after the Learning process.**

## Models:

### Machine-Learning , Neural-Network:

- 1) Deep Convolutional Neural Network with language Embeddings Representation.
- 2) Siamese Neural networks with language Embeddings Representation.

- \*embedding\_model
- \*embedding\_model2
- \*embedding\_lstm\_model\_manhattan\_dist (Siamese)
- \*embedding\_lstm\_model (Siamese)

## Evaluation:

**K-Fold.**

## Data-Sets:

Dataset 1: Recipe description vs Recipe review

Dataset2: Covid19-questions vs Covid19-answer

Dataset 3: Questions-Pairs. (question1 vs question2)

## Dataset 1: Recipe description vs Recipe review

About: Recipe dataset, very nice dataset to work with, a lot of options.

Viewer	Text																										
<div><div>JSON</div><div><div>0</div><div><div>name : "arriba baked winter squash mexican style"</div><div>id : "137739"</div><div>minutes : "55"</div><div>contributor_id : "47892"</div><div>submitted : "09/16/2005"</div><div>tags : ["60-minutes-or-less", "time-to-make", "course", "main-ingredient", "cuisine", "preparation", "occasion", "north-american", "side-dishes", "veg"</div><div>nutrition : "[51.5, 0.0, 13.0, 0.0, 2.0, 0.0, 4.0]"</div><div>n_steps : "11"</div><div>steps : ["make a choice and proceed with recipe", "depending on size of squash , cut into half or fourths", "remove seeds", "for spicy squash ,"</div><div>description : "autumn is my favorite time of year to cook! this recipe can be prepared either spicy or sweet, your choice! two of my posted m"</div><div>ingredients : ["winter squash", "mexican seasoning", "mixed spice", "honey", "butter", "olive oil", "salt"]"</div><div>n_ingredients : "7"</div></div></div><div><div>1</div><div><div>name : "a bit different breakfast pizza"</div><div>id : "31490"</div><div>minutes : "30"</div><div>contributor_id : "26278"</div><div>submitted : "06/17/2002"</div><div>tags : ["30-minutes-or-less", "time-to-make", "course", "main-ingredient", "cuisine", "preparation", "occasion", "north-american", "breakfast", "main-</div><div>nutrition : "[173.4, 18.0, 0.0, 17.0, 22.0, 35.0, 1.0]"</div><div>n_steps : "9"</div><div>steps : ["preheat oven to 425 degrees F", "press dough into the bottom and sides of a 12 inch pizza pan", "bake for 5 minutes until set but not"</div><div>description : "this recipe calls for the crust to be prebaked a bit before adding ingredients. feel free to change sausage to ham or bacon. this"</div><div>ingredients : ["prepared pizza crust", "sausage patty", "eggs", "milk", "salt and pepper", "cheese"]"</div><div>n_ingredients : "6"</div></div></div><div><div>2</div><div><div>name : "all in the kitchen chili"</div><div>id : "112140"</div><div>minutes : "130"</div><div>contributor_id : "196586"</div><div>submitted : "02/25/2005"</div><div>tags : ["time-to-make", "course", "preparation", "main-dish", "chili", "crock-pot-slow-cooker", "dietary", "equipment", "4-hours-or-less"]"</div><div>nutrition : "[269.8, 22.0, 32.0, 48.0, 39.0, 27.0, 5.0]"</div><div>n_steps : "6"</div><div>steps : ["brown ground beef in large pot", "add chopped onions to ground beef when almost brown and sautee until wilted", "add all other ing"</div><div>description : "this modified version of 'mom's' chili was a hit at our 2004 christmas party. we made an extra large pot to have some left to fre"</div><div>ingredients : ["ground beef", "yellow onions", "diced tomatoes", "tomato paste", "tomato soup", "rotel tomatoes", "kidney beans", "water", "chili poi"</div><div>n_ingredients : "13"</div></div></div></div>	<table><tr><th>Name</th><th>Value</th></tr><tr><td>contributor_id</td><td>"47892"</td></tr><tr><td>description</td><td>"autumn is my favorit..."</td></tr><tr><td>id</td><td>"137739"</td></tr><tr><td>ingredients</td><td>"['winter squash', 'me..."</td></tr><tr><td>minutes</td><td>"55"</td></tr><tr><td>name</td><td>"arriba baked winter s..."</td></tr><tr><td>nutrition</td><td>"[51.5, 0.0, 13.0, 0.0, ..."</td></tr><tr><td>n_ingredients</td><td>"7"</td></tr><tr><td>n_steps</td><td>"11"</td></tr><tr><td>steps</td><td>"['make a choice and ..."</td></tr><tr><td>submitted</td><td>"09/16/2005"</td></tr><tr><td>tags</td><td>"['60-minutes-or-less', ..."</td></tr></table>	Name	Value	contributor_id	"47892"	description	"autumn is my favorit..."	id	"137739"	ingredients	"['winter squash', 'me..."	minutes	"55"	name	"arriba baked winter s..."	nutrition	"[51.5, 0.0, 13.0, 0.0, ..."	n_ingredients	"7"	n_steps	"11"	steps	"['make a choice and ..."	submitted	"09/16/2005"	tags	"['60-minutes-or-less', ..."
Name	Value																										
contributor_id	"47892"																										
description	"autumn is my favorit..."																										
id	"137739"																										
ingredients	"['winter squash', 'me..."																										
minutes	"55"																										
name	"arriba baked winter s..."																										
nutrition	"[51.5, 0.0, 13.0, 0.0, ..."																										
n_ingredients	"7"																										
n_steps	"11"																										
steps	"['make a choice and ..."																										
submitted	"09/16/2005"																										
tags	"['60-minutes-or-less', ..."																										

Original Files: RAW\_recipes.csv and RAW\_interactions.csv

Preprocess: Remove unwanted chars.

Export: pairs file **recpies\_DS\_1.csv**

Out[8]:

	description	review
0	autumn is my favorite time of year to cook! th...	I used an acorn squash and recipe#137681 Swee...
1	autumn is my favorite time of year to cook! th...	This was a nice change. I used butternut squas...
2	autumn is my favorite time of year to cook! th...	Excellent recipe! I used butternut squash and ...
3	this recipe calls for the crust to be prebaked...	Have not tried this, but it sounds delicious. ...
4	this recipe calls for the crust to be prebaked...	This recipe was wonderful. Instead of using t...

## Dataset2: General questions vs right or wrong answer

Original Files: general.csv

View:

ViewerText

JSON

0

- question\_id : "18276"
- title : "Questions on the effectiveness and safety of weighted pistol squats"
- question : "I recently stopped doing front squats due to bicep tendonitis in my left arm and I peaked at 215 for 5 reps. So now since I am un
- answer\_id : "18290"
- answer : "DOMS is not directly indicative of how effective a workout is, so because you don't feel them doesn't mean its not an effective wor
- answer\_type : "Accepted"
- wrong\_answer : "I'm a beginner (40 years old) who started running the 5k three times per week. It took about three or four weeks but I s
- wrong\_answer\_type : "Random"
- site : "fitness"
- group : "general"

1

- question\_id : ""
- title : ""
- question : ""
- answer\_id : ""
- answer : "It depend on how you define "knee". Since pistol squats will be done with less weight on a per knee basis the joint itself will under
- answer\_type : ""
- wrong\_answer : ""
- wrong\_answer\_type : ""
- site : ""
- group : ""

2

- question\_id : ""
- title : ""
- question : "I have three questions"
- answer\_id : ""
- answer : "Given that pistol squats require significantly more work from your assistance muscles to maintain balance its very hard to estimat
- answer\_type : ""
- wrong\_answer : ""
- wrong\_answer\_type : ""
- site : ""
- group : ""

Name	Value
answer	"Given that pistol squa...
answer_id	""
answer_type	""
group	""
question	"I have three questions"
question_id	""
site	""
title	""
wrong_answer	""
wrong_answer_type	""

Preprocess:

1. Remove unwanted words that does not fit the regression rule: [A-Za-z0-9] and between 2-50 chars.
2. Remove stop words using the stopwords.txt file
3. Lemmatization
4. Replace capital letters

Export: file **precovid\_data.csv**

In [7]: data.head()

Unnamed: 0	question	answer	wrong_answer
0	I recently stopped doing front squats due to b...	DOMS is not directly indicative of how effec...	I'm a beginner (40 years old) who started run...
1	I have a lot of fat lose. I weigh 210 pounds, ...	Yes, its theoretically possible. Are you going...	Welcome to the Fitness SE!n/nMaking your own ...
2	After going through this link I am confused ab...	It appears that this article is aimed more at ...	There is a myth that exercise increases the se...

## Dataset 3: Questions-Pairs. (question1 vs question2)

Original Files: questions.csv

View:

```
In [6]: data.head()
```

Out[6]:

	Unnamed: 0	question1	question2	is_duplicate
0	0	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0
1	1	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0
2	2	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0

### Preprocess 1:

1. Remove unwanted words that does not fit the regression rule: [A-Za-z0-9] and between 2-50 chars.
2. Remove stop words using the stopwors.txt file
3. Lemmatization
4. Replace capital letters

Export: file **p\_c\_questions\_data.csv**

### Preprocess 2:

1. Remove unwanted words that does not fit the regression rule: [A-Za-z0-9] and between 1-50 chars.
2. Remove stop words using the stopwors.txt file (updated the stopwors.txt)
3. Lemmatization
4. Replace capital letters

Export: file **p\_c\_questions\_data2.csv**

### Preprocess3:

1. Remove unwanted words that does not fit the regression rule: [A-Za-z0-9] and between 1-50 chars.
2. Lemmatization
3. Replace capital letters

Export: file `p_c_questions_data3.csv`

## Project Files description

FinalDataGenerator1.py :

Knows how to work with this type of scheme :

Out[8]:

	description	review
0	autumn is my favorite time of year to cook! th...	I used an acorn squash and recipe#137681 Swee...
1	autumn is my favorite time of year to cook! th...	This was a nice change. I used butternut squas...
2	autumn is my favorite time of year to cook! th...	Excellent recipe! I used butternut squash and ...
3	this recipe calls for the crust to be prebaked...	Have not tried this, but it sounds delicious. ...
4	this recipe calls for the crust to be prebaked...	This recipe was wonderful. Instead of using t...

Only positive samples.

At the main function:

You can generate batches with your own parameters.

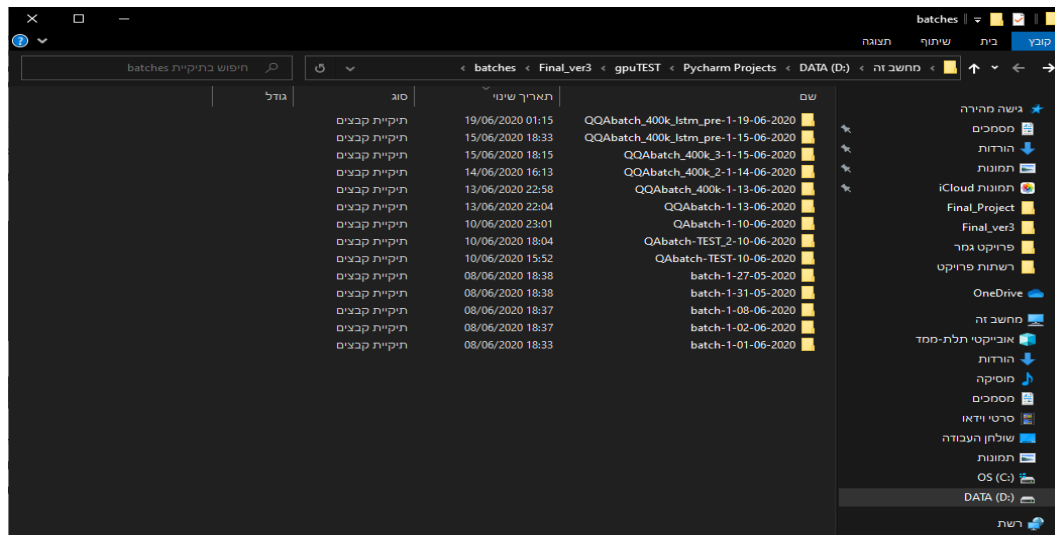
```
batch1 = next(self.generate_batch(1000,1)) # 2000
batch1 = self.arrange_batch(batch1,self.batch_length)
self.save_batch(batch1,1000,1)
batch2 = next(self.generate_batch(self.num_of_rows,1))
batch2 = self.arrange_batch(batch2,self.batch_length)
self.save_batch(batch2,self.num_of_rows,1)

def generate_batch(self, n_positive=None,
negative_ratio=1):
```

you can change the negative ratio as you like. The default is 1 to 1.

This function yields the batch.

The save\_batch function saves the batch to the "batches" folder and creates there your batches.



FinalDataGenerator2.py :

Knows how to work with this type of scheme :

In [7]: `data.head()`

Out[7]:

	Unnamed: 0	question	answer	wrong_answer
0	0	I recently stopped doing front squats due to b...	\nDOMS is not directly indicative of how effec...	I'm a beginner (40 years old) who started run...
1	1	I have a lot of fat lose. I weigh 210 pounds, ...	Yes, its theoretically possible. Are you going...	Welcome to the Fitness SE!\n\nMaking your own ...
2	2	After going through this link I am confused ab...	It appears that this article is aimed more at ...	There is a myth that exercise increases the se...

With positive and negative samples at the same file.(without labels)

At the main function:

```
def main(self):
    #self.load_files(500)
    self.pre_process()
    self.load_files(130000)
    self.data_preperation()

    batch2 = next(self.generate_batch(130000,1))
    batch2 = self.arrange_batch(batch2,self.batch_length)
    self.save_batch(batch2,self.num_of_rows,1)
```

Notice the pre\_process() Function , It takes the data and run it through the preprocess pipe line, you can remove or add functions as you like.

It exports the data to this kind of scheme:

```
In [7]: data.head()
```

```
Out[7]:
```

Unnamed: 0	question	answer	wrong_answer	clean_question	clean_answer	clean_wrong_answer	p_c_question	p_c_answer	p_c_wrong_answer
0	0	I recently stopped doing front squats due to b...	inDOMS is not directly indicative of how effec...	I'm a beginner (40 years old) who started runn...	recently stopped doing front squats due to bic...	DOMS is not directly indicative of how effecti...	beginner 40 years old who started running the ...	[recently, stop, squat, bicep, tendonitis, lea...	[doms, directly, indicative, effective, workout...
1	1	I have a lot of fat lose. I weigh 210 pounds, ...	Yes, its theoretically possible. Are you going...	Welcome to the Fitness SE!n/nMaking your own ...	have lot of fat lose weigh 210 pounds and was ...	Yes its theoretically possible Are you going t...	Welcome to the Fitness SE Making your own prog...	[lot, fat, lose, weigh, pound, hope, start, lo...	[yes, theoretically, possible, able, decide, s...
2	2	After going through this link I am confused ab...	It appears that this article is aimed more at ...	There is a myth that exercise increases the se...	After going through this link am confused about...	It appears that this article is aimed more at ...	There is myth that exercise increases the secr...	[link, confuse, believe, link, claim, true, af...	[appear, article, aim, condition, activity, ru...

Generate batch works a little bit different because of the data set format but yields and saves the results same as the first one.

FinalDataGenerator3.py :

Knows how to work with this type of scheme :

```
In [6]: data.head()
```

```
Out[6]:
```

Unnamed: 0	question1	question2	is_duplicate
0	0	What is the step by step guide to invest in sh...	0
1	1	What is the story of Kohinoor (Koh-i-Noor) Dia...	0
2	2	How can I increase the speed of my internet co...	0

With positive and negative samples at the same file already have labels to each pair.

The preprocess part same as: FinalDataGenerator2.py

```
In [6]: data.head()
```

```
Out[6]:
```

Unnamed: 0	question1	question2	is_duplicate	clean_question1	clean_question2	p_c_question1	p_c_question2
0	0	What is the step by step guide to invest in sh...	0	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	[step, step, guide, invest, share, market, india]	[step, step, guide, invest, share, market]
1	1	What is the story of Kohinoor (Koh-i-Noor) Dia...	0	What is the story of Kohinoor Koh-i-Noor Diamond	What would happen if the Indian government sto...	[story, kohinoor, koh, noor, diamond]	[happen, indian, government, steal, kohinoor, ...]
2	2	How can I increase the speed of my internet co...	0	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	[increase, speed, internet, connection, vpn]	[internet, speed, increase, hack, dns]

Generate batch works a little bit different because of the data set format but yields and saves the results same as above.

FinalLoadBatch.py

loading the saved batch from the batches folder.

```
class LoadBatch:
    def __init__(self, dir_name):
```

FinalModels.py

Run:

```
a = Models("embedding_lstm model2", "QQAbatch 400k lstm pre-1-2-06-2020")
model = a.model
print(model.summary())
model = a.start()
```

The first parameter is the model you want to run.

The second parameter is the name of the batch file at the "batches" directory. Loads the data using "LoadBatch" Object.

At the "start" function you can set the number of K at k-fold evaluation.

FinalGUI.py



```
def loadData(self):
    # self.data = LoadBatch('batch-1-07-05-2020')
    self.data = LoadBatch('QQAbatch_400k_lstm_pre-1-19-06-2020')
    self.tokenizer = self.data.tokenizer
    self.max_len_description = self.data.left_vector_max_size
    self.max_len_review = self.data.right_vector_max_size
    self.max_words = self.data.vocab_size
    self.df = self.data.convert_to_dataframe()
    # self.model = load_model('embedding_model2_07-05-2020.h5')
    self.model = load_model('400k_model_embedding_TEST-LIAD_3_QQAbatch-1-20-06-2020.h5')
    # self.model = load_model('400k_model_embedding_TEST-LIAD_3_QQAbatch-1-21-06-2020.h5')
```

First, you need to load the batch you used to train the model.

Second, you need to load the Neural Network model you have trained as h5 file.