

## What has been done?

- Run full w2v on AWS, on strongest machine available → expensive
- Testing Azure Notebook → storage limitation
- Testing on my machine (PC) → ram limitation
- Made an iteration/generator methode to do line by line on dataset
- Made a wikipedia dataset splitter with balise preservation
- Pre-trained models are not retrainable, not full model
- Dictionary unexpected: - Lemmatization missing → Only Keyed Vectors  
→ Dog/NN, Do/VB, ...  
"library not installed first time run, "pattern" "
- 1 CPU Core not working
- Memory allocation error on splitter also

## short term

- starting new sprint on - ANN chatbot
  - parallel algorithm
  - protocole to evaluate proactive

Get

- full w2v wikipedia
- meanwhile use smaller dataset
- Creation and test of existing chatbot solution
- Fix memory error on splitter with a buffer for the generator

## Last time short term

- Combine with ANN
- Document in report Gensim life
- Compare with: FastText, word2Vec, Glove

## overall progress

- fixing the memory problem with gensim by splitting datasets by pages and use a generator by line
- Trying to use cloud based machines, but expensive and no results yet.
- Word2Vec operations working on partial model by dictionary is not as expected

## Questions:

- Problématique de recentrage
- Tester les translations  $\rightarrow$  avec un set de phrases  $\rightarrow$  20 ans
- Case comprehension Word2Vec  $\rightarrow$

## Translation

- $\rightarrow$  perturbation mini
- $\rightarrow$  le chat est une mammifère
  - $\rightarrow$  Plus proche voisin
    - $\rightarrow$  vecteur appliquer au reste de la phrase
    - $\rightarrow$  vecteur de translation
- $\rightarrow$  translation dirigé  $\rightarrow$  shift
- $\rightarrow$  translation random
- $\rightarrow$  Word2Vec
  - $\rightarrow$  A la base prétraitement des mots
  - $\rightarrow$  Playing with word2vec
- $\rightarrow$  Dimensions  $\rightarrow$  concepts
  - $\rightarrow$  Jump on 1 dimension
  - $\rightarrow$  Multi dimensions
- $\rightarrow$  Synonyme
  - $\rightarrow$  point chat
    - $\rightarrow$  synonyme de chat
    - $\rightarrow$  quelles dimensions les plus impactées
    - $\rightarrow$  dériver dans la direction
- But comprendre l'espace Word2Vec
- $\rightarrow$  Word2Sequence

- Chapitre sur le biais

-> mettre en évidence

Doctor - man + woman = Nurse

-> Facebook, Google, -> Biais dans recommandation  
publicité

-> poussé vers des études