

When was this sentence written?

bertrand.gaiffe@atilf.fr

may 25th, 2021



The Context

Some of us are interested in older states of the french language :

- DMF (dictionnaire of middle french)
- Frantext...
- a whole team in the lab (historical linguistics of french and romance languages)

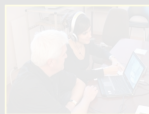
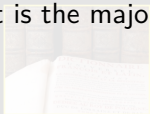
Therefore, why not a BERT model for (in particular) middle french?

- a group of people (B. Crabbé, M. Constant, L. Grobol, P. Ortiz, T. Mickus...)
- a corpus of the form (DATE SENTENCE)+

The corpus

- Frantext (texts before 1920)
- charters
- Ambroise Paré, Champlain, etc.
- OEF (Ovide en Français livre 1)
- Ovide moralisé en prose

Frantext is the major part....



Why not guessing the date?

Because we have this corpus, why not building a date guesser ?

So the problem is to guess the date when a “sentence” was written.

Examples :

- - Attendez donc, reprit le conteur ; voici le plus beau et le plus terrible de l’histoire que je vous avais promise : Coupe-en-Deux avait tombé par terre comme un plomb ; il était si soûl, si soûl, qu’il ne remuait pas plus qu’une bûche..
- - Renart, dist Nobles, bien as dit, ja de ce n’i avras desdit: or diras, nos t’escoterons; se tu diz bien, nos nos tairons.
- Helene Vous estes fou, Mendoce. Dom Diegue Vous estes folle, Helene, avecque vostre nopce.

Why not guessing the date?

Because we have this corpus, why not building a date guesser ?

So the problem is to guess the date when a “sentence” was written.

Examples :

- - Attendez donc, reprit le conteur ; voici le plus beau et le plus terrible de l’histoire que je vous avais promise : Coupe-en-Deux avait tombé par terre comme un plomb ; il était si soûl, si soûl, qu’il ne remuait pas plus qu’une bûche.. (1843)
- - Renart, dist Nobles, bien as dit, ja de ce n’i avras desdit: or diras, nos t’escoterons; se tu diz bien, nos nos tairons.
- Helene Vous estes fou, Mendoce. Dom Diegue Vous estes folle, Helene, avecque vostre nopce.

Why not guessing the date?

Because we have this corpus, why not building a date guesser ?

So the problem is to guess the date when a “sentence” was written.

Examples :

- - Attendez donc, reprit le conteur ; voici le plus beau et le plus terrible de l’histoire que je vous avais promise : Coupe-en-Deux avait tombé par terre comme un plomb ; il était si soûl, si soûl, qu’il ne remuait pas plus qu’une bûche.. (1843)
- - Renart, dist Nobles, bien as dit, ja de ce n’i avras desdit: or diras, nos t’escoterons; se tu diz bien, nos nos tairons. (1190)
- Helene Vous estes fou, Mendoce. Dom Diegue Vous estes folle, Helene, avecque vostre nopce.

Why not guessing the date?

Because we have this corpus, why not building a date guesser ?

So the problem is to guess the date when a “sentence” was written.

Examples :

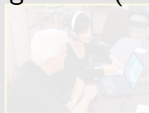
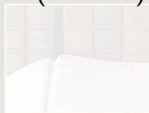
- - Attendez donc, reprit le conteur ; voici le plus beau et le plus terrible de l’histoire que je vous avais promise : Coupe-en-Deux avait tombé par terre comme un plomb ; il était si soûl, si soûl, qu’il ne remuait pas plus qu’une bûche.. (1843)
- - Renart, dist Nobles, bien as dit, ja de ce n’i avras desdit: or diras, nos t’escoterons; se tu diz bien, nos nos tairons. (1190)
- Helene Vous estes fou, Mendoce. Dom Diegue Vous estes folle, Helene, avecque vostre nopce. (1650)

Characteristics of the corpus

A large part of the corpus is from Frantext. Therefore, a lot of XIXth century sentences.

Most basic baselines :

- most frequent year (on train) : 1910 (average error (on test) = 107.75)
- mean (in train) : 1805 (average error (on test) = 96.44)

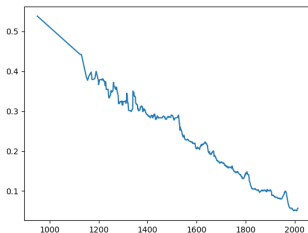


An experiment about trigrams

Total variation distance (date 1, date 2):

$$\sum_{tri} |P(tri/d1) - P(tri/d2)|$$

(because we pondered whether fine tuning an existing BERT...)



tvd(d, 2000) smoothed on 20 years

Trigrams and Bayes rule

- Compute $P(\text{Trigramme}/\text{Date})$ for all trigrams
- and also $P(\text{trigramme})$ and $P(\text{date})$
- then $P(\text{date}/\text{trigramme}) = \frac{P(\text{trigramme}/\text{date}) \times P(\text{date})}{P(\text{trigramme})}$

and then, we do something absolutely wrong (because no independance), that is :

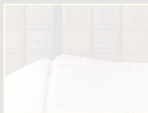
$$P(\text{date}/t_1, t_2, \dots, t_n) = \frac{P(t_1/\text{date}) \times P(t_2/\text{date}) \dots \times P(\text{date})}{P(t_1) \times P(t_2) \dots}$$

Entirelly false, but we are interested in $\arg\text{Max}$ for the dates anyway....

Results with trigrams

- We compute (count) the probas on TRAIN and we test on TEST.
- Result : mean error = 31 years ! (31.47)

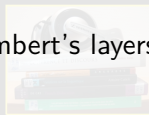
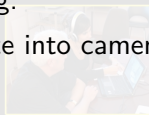
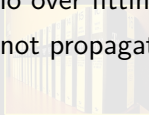
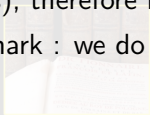
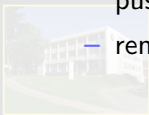
Remember : mean date (1805) gives an average error of 96 years



Using BERT for datation

Experiment :

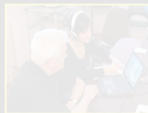
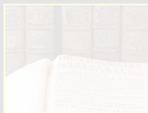
- Take some of the layers, and add a linear classifier on them, take the mean, then relu, tanh and normalize between 900 and 1920
- Remark, we consider dates as numbers (not categories)
- we do not even learn a full epoch (just around 5% of the corpus), therefore no over fitting.
- remark : we do not propagate into camembert's layers



results

- layer n° 12 \longrightarrow 68 years
- layer n° 1 \longrightarrow 98 years

These results correspond to about 1/20th of an epoch. Therefore, no overfitting !

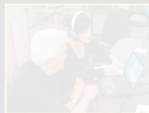


Is it a wordpieces problem ?

The Camembert model relies on wordpieces (not trigrams).

Could that be the problem ?

Well, same strategy as trigrams but with wordpieces gives 35 years error in average. (as compared to 31).



Very partial conclusions

- should compute not only mean error but variance (and Bert model should be better on this)
- should smooth the $P(tri/date)$ in the trigram approach.
- test with several layers (see whether different informations) (first test : layer 12 + layer 9 not better than layer 12 alone.

