# PARTICIPEZ À UNE COMPÉTITION KAGGLE !

OPENCLASSROOMS - INGÉNIEUR MACHINE LEARNING

THIBAUD GROSJEAN - MARS 2022

# OBJECTIFS DU PROJET

- Participer à une compétition *Kaggle ("The House of Data Science")*

- Utiliser les ressources partagées par la communauté

- Partager un élément intéressant avec la communauté

- Expliciter le modèle selectionné

# COMPÉTITION

- **NBME – Score Clinical Patient Notes**
  - *United States Medical Licensing Examination® (USMLE®)*
  - Examen de compétences cliniques
  - Simulation de consultation
  - Historique du patient

# VERSIONNAGE

# ENVIRONNEMENT

- *Kaggle Notebooks*

- *Internet*

- *Parallelisation*

## OBJECTIF

- *Natural Langage Processing (NLP)*
- *Named-entity Recognition (NER)*
- Evaluation (*score F1 micro-moyenné*)

# TRANFORMEURS

- *BERT*

- *RoBERTa*

- Modèles *auto-supervisés*

# DONNÉES

- *Patient Notes* (42146)

- *Features* (143)

- Création d'un fichier plat
  (14300 samples)

# EMBEDDINGS

- Vecteurs representant numériquement les caractéristiques du *document*
- *Input ids* (tokens)
- *Attention masks*
- *Labels* (cibles)

# PIPELINE

- Simple en apparence…
- Mais complexe
- Le *Tokenizer* permet d'inverser les *logits* du modèle

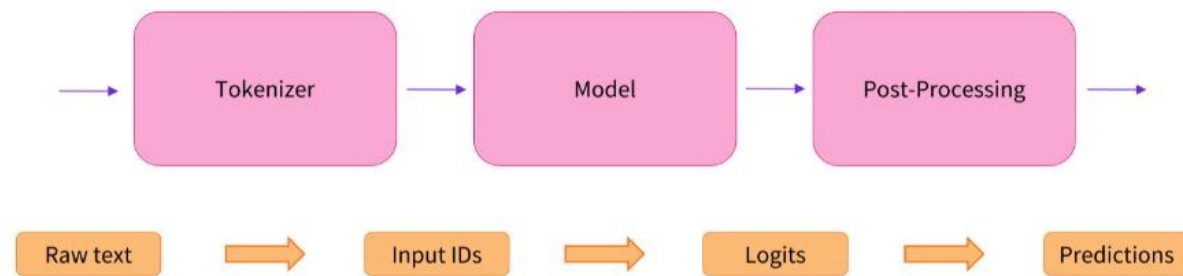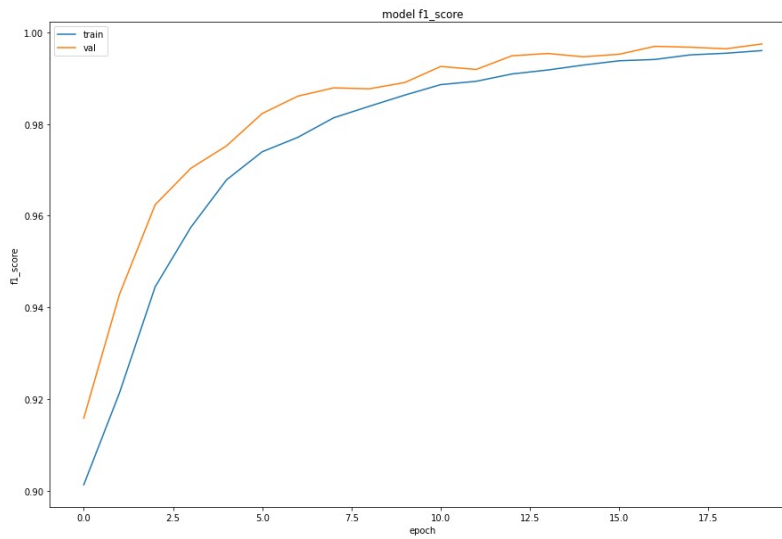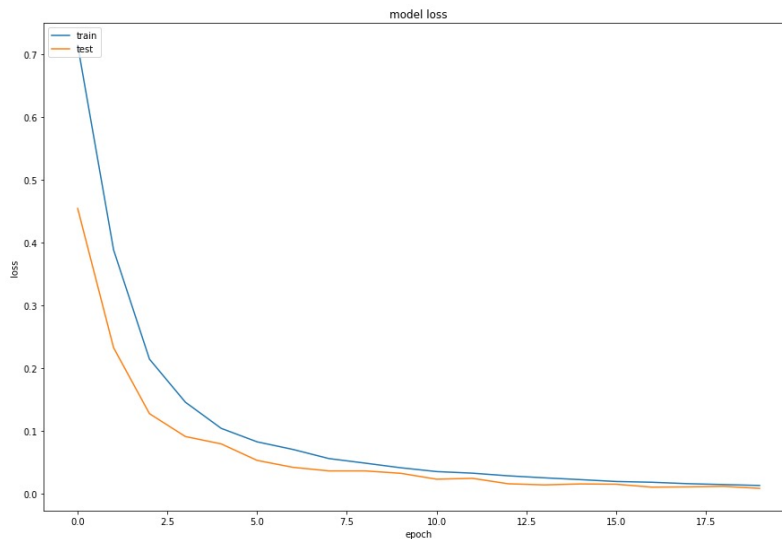| Comparison | BERT October 11, 2018 | RoBERTa July 26, 2019 | DistilBERT October 2, 2019 | ALBERT September 26, 2019 |
|---|---|---|---|---|
| Parameters | **Base:** 110M **Large:** 340M | **Base:** 125 **Large:** 355 | **Base:** 66 | **Base:** 12M **Large:** 18M |
| Layers / Hidden Dimensions / Self-Attention Heads | **Base:** 12 / 768 / 12 **Large:** 24 / 1024 / 16 | **Base:** 12 / 768 / 12 **Large:** 24 / 1024 / 16 | **Base:** 6 / 768 / 12 | **Base:** 12 / 768 / 12 **Large:** 24 / 1024 / 16 |
| Training Time | **Base:** 8 x V100 x 12d **Large:** 280 x V100 x 1d | 1024 x V100 x 1 day (4-5x more than BERT) | **Base:** 8 x V100 x 3.5d (4 times less than BERT) | [not given] **Large:** 1.7x faster |
| Performance | Outperforming SOTA in Oct 2018 | 88.5 on GLUE | 97% of BERT-base's performance on GLUE | 89.4 on GLUE |
| Pre-Training Data | BooksCorpus + English Wikipedia = 16 GB | BERT + CCNews + OpenWebText + Stories = 160 GB | BooksCorpus + English Wikipedia = 16 GB | BooksCorpus + English Wikipedia = 16 GB |
| Method | Bidirectional Transformer, MLM & NSP | BERT without NSP, Using Dynamic Masking | BERT Distillation | BERT with reduced parameters & SOP (not NSP) |

# AMÉLIORATION DU SCORE

- Parallélisation

- Split des données

- Visualisation

- Implémentation d'un modèle plus performant *(RoBERTa)*

- Optimisation du nombre d'*epochs*

- *Model Checkpoint*

# ENTRAINEMENT & OPTIMISATION

- Le modèle *merge* correctement

# AMÉLIORATION DU SCORE

# AMÉLIORATION DE LA RESSOURCE EXISTANTE

- <u>Création d'une ressource dictatisée</u>

- Ajout de commentaires & de ressources

- Simplification du process *Kaggle notebooks*

- Amélioration du style

- <u>Amélioration de la méthodologie en *Data Science*</u>

# CONCLUSION & PISTES D'AMÉLIORATION

- Prise de main de *Kaggle Notebooks* & de la *Parallelisation*

- Découverte des *Transformers* & de *Hugging Face*

- Utilisation des ressources communautaires

- Amélioration du score

- Partage de la ressource créée

# ÉCHANGE & QUESTIONS

Merci de votre attention !