



Apprentissage Profond

Génération de proverbes

Groupe L34-E

Élèves :

THEVENET Louis

LEBOBE Timothé

TENE Zacharie

SABLAYROLLES Guillaume

31 Mars 2025

Table des matières

1. Génération de proverbes en anglais	3
2. Consitution de la base de données	3
2.1. Acquisition des données	3
2.2. Script de chargement des données	3
3. Création et entraînement du modèle	4
3.1. Création du modèle	4
3.2. Entraînement	5
3.3. Génération	5
4. Analyse des résultats	6
4.1. Génération de proverbes sans thème	6
4.2. Génération avec thème	11
5. Conclusion	11

1. Génération de proverbes en anglais

Nous avons choisi de créer un modèle de génération de proverbes. La base de donnée est trouvable dans le dossier `raw_data/` à la racine de ce [dépôt GitHub](#)

Voici un exemple de proverbes de notre base d'entraînement :

- He that brings good news, knocks hard.
- Anger and haste hinder good counsel.
- Big thunder, little rain.
- Romeo must die in order to save the love.
- The point is plain as a pike staff.

2. Consitution de la base de données

2.1. Acquisition des données

Puisqu'il est plus simple de trouver des données en langue anglaise, nous avons choisi de nous limiter à cette langue et avons utilisé des scripts Python de scrapping pour récolter des données sur différents sites internet.

Nous avons 3200 proverbes originaux anglais et 35000 proverbes en incluant des proverbes traduits d'autres langues et prévoyons de tester le modèle sur ces deux bases.

Nous également constitué une base de données de proverbes annotés avec leurs thèmes.

Voici quelques exemples :

```
1  {
2    "topics": ["Advantage"],
3    "proverb": "When a rich man caresses a poor man, he's going to take advantage
4    of him."
5  }
6  {
7    "topics": ["Advantage"],
8    "proverb": "Every advantage has its disadvantage."
9  }
10 {
11    "topics": ["Advantage", "Vain"],
12    "proverb": "He is wise in vain who does not use his wisdom for his own
    advantage."
```

Liste 1. – Extrait des proverbes annotés

2.2. Script de chargement des données

Nous avons réalisé un script de téléchargement et traitement des données. Un exemple d'utilisation est donné dans le fichier `main.ipynb`, il suffit d'appeler la fonction `make_dataset.load_data()` qui renvoie les proverbes classés par sources.

```

1  proverbs_db.txt: 34142 proverbs
2  proverbs_db_only_english.txt: 2208 proverbs
3  proverbs_digest.txt: 1000 proverbs
4  Total: 37350 proverbs
5  Total length: 1853527 characters
6
7
8  Examples of proverbs:
9  Money's for buying and a horse is for riding.
10 A set of white teeth does not indicate a pure heart.
11 Time discloses the truth.
12 The earth has ears, the wind has a voice.
13 The fox will catch you with cunning, and the wolf with courage.

```

Liste 2. – Proverbes non annotés

De même, la fonction `make_dataset.load_data_with_topics()` renvoie un dataset de données annotées.

3. Création et entraînement du modèle

3.1. Création du modèle

Nous allons faire du fine-tuning à partir du modèle [facebook/opt-125](#) qui est un modèle type GPT-3 basé sur l'architecture Transformers.

Nous testerons également de partir du modèle [TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T](#) qui est plus récent et contient près de dix fois plus de paramètres.

On charge le modèle et son Tokenizer à l'aide des fonctions `AutoModelForCausalLM.from_pretrained()` et `AutoTokenizer.from_pretrained()` de la librairie `transformers`.

On décide des sources de proverbes que l'on va utiliser, puis on les fusionne en une seule liste.

```

1  selected_proverbs_groups = [
2      "proverbs_db.txt",
3      "proverbs_digest.txt"
4  ]
5
6  proverbs = []
7  for group in selected_proverbs_groups:
8      proverbs.extend(all_proverbs[group])

```

Liste 3. – Fusion des listes sélectionnées

On utilise ensuite la librairie `datasets` pour préparer ces données à l'entraînement. La tokenisation du dataset consiste à calculer la taille du proverbe le plus long et ajouter du padding aux autres pour uniformiser les tailles.

On utilise ensuite la librairie `peft` afin d'utiliser la technique LoRA (Low-Rank Adaptation). Ainsi, on ajoute un petit nombre de nouveaux paramètres entraînaables au modèle afin de l'adapter à la nouvelle tâche.

```

1 LoraConfig(
2     r=8,
3     lora_alpha=16,
4     target_modules=["q_proj", "v_proj"],
5     lora_dropout=0.05,
6     bias="none",
7     task_type="CAUSAL_LM"
8 )

```

Liste 4. – Configuration LoRA

La méthode `get_peft_model` permet d'obtenir un nouveau modèle à partir de cette configuration de notre modèle initial.

3.2. Entraînement

A l'aide de la librairie `transformers`, on définit les paramètres d'entraînement:

```

1 TrainingArguments(
2     output_dir="./results",
3     per_device_train_batch_size=8,
4     per_device_eval_batch_size=8,
5     num_train_epochs=1,
6     logging_dir='./logs',
7     logging_steps=10,
8     eval_strategy="no",
9     save_strategy="epoch",
10    report_to="none"
11 )

```

Liste 5. – Paramètres d'entraînement du modèle

Finalement, on met en commun notre modèle, nos paramètres d'entraînement, notre dataset tokenisé et notre tokenizer via la classe `Trainer` et on peut lancer l'entraînement avec la méthode `train()`.

3.3. Génération

Après entraînement du modèle, on crée un pipeline de génération.

```

1 generator = pipeline("text-generation", model=model, tokenizer=tokenizer)

```

Liste 6. – Pipeline de génération

On peut maintenant utiliser le modèle pour terminer un début de proverbes:

```

1 prompt = "A"
2 results = generator(prompt, max_length=max_length, num_return_sequences=3,
3 do_sample=True, temperature=0.7)
4 for i, result in enumerate(results):
5     print(result['generated_text'])

```

Liste 7. – Génération de proverbes

Quelques proverbes obtenus:

- A nice dog is a dog.
- A man who keeps his family safe can never be found.
- A child that has no teeth is a coward.

4. Analyse des résultats

Les données de sortie étant des proverbes, elles sont difficilement évaluables et comparables, et sont sujettes à l'appréciation humaine.

Ainsi nous donnerons quelques résultats non filtrés à titre d'exemple, issus du fine tuning de différents modèles que nous commenterons. Puis dans une seconde partie, nous détaillerons une expérience que nous avons réalisé.

4.1. Génération de proverbes sans thème

4.1.1. Entraînement sur les proverbes anglais

On utilise ici les deux sources :

- `proverbs_db_only_english.txt`
- `proverbs_digest.txt`

Pour rappel, cette base de données représente 3208 proverbes originaux en anglais.

Nous allons classer les générations suivant 4 interprétations : *Intéressants*, *Etranges*, *Pas un proverbe* et *Incohérents*. Cela représente comment nous percevons ce qu'est un proverbe car, vous verrez en suivant, certaines générations donnent des sortes de vérités générales (que nous pointerons) qui ne ressemblent pas à ce qu'est un proverbe pour nous.

Cependant ces vérités générales se rapprochent de certains proverbes de notre base de données (cf Conclusion).

4.1.1.1. Modèle de départ `facebook/opt-125m`

Proverbe	Commentaire
A man's only dead when he eats his own.	Intéressant
A man can make a woman forget her brother's death.	Pas un proverbe
A man's heart is full of gold.	Pas un proverbe
A good thing is a good thing.	Pas un proverbe
A woman's dream is to be a widow.	Etrange
A man dies in the wind, a horse dies in the wind	Pas un proverbe
A good example of a good example of a good example of a bad example of a bad example of a bad example.	Incohérent
A few days before you start a new one, you will be remembered for a long time.	Incohérent
A little of light can be a good thing	Pas un proverbe

Tableau 1. – En partant de « A »

Proverbe	Commentaire
Some people have a hard time keeping a family man.	Intéressant
Some people are good, some people are good.	Incohérent
Some people are lucky to live in a tree.	Intéressant
Some of them are to be proud of their first friend.	Intéressant
Some times the best is when you are there to be with.	Incohérent
Some people have a disease.	Etrange
Some people are better than others.	Intéressant
Some of us are better than others	Intéressant
Some people are too silly to forget.	Intéressant
Some of those is nothing compared to the sum of them	Intéressant

Tableau 2. – En partant de « Some »

Intéressants 9

Etranges 2

Pas un proverbe 5

Incohérents 4

On constate que certains résultats sont incohérents ou étranges. Certains proverbes ressemblent plus à des vérités générales et on trouve quelques proverbes intéressants.

4.1.1.2. Modèle de départ **TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T**

Les données sélectionnées ne sont pas assez importantes pour obtenir des résultats satisfaisants avec ce modèle.

Proverbe	Commentaire
This is the story of an old woman.	Pas un proverbe
This is the way of the world.	Pas un proverbe
This is the best.	Pas un proverbe
This is the way to get your money's worth.	Pas un proverbe
This is the day.	Pas un proverbe
This is what you call a new dress.	Pas un proverbe
This is a story about a little girl.	Pas un proverbe
This is the day the Lord hath spoken. It is the day that the Lord hath spoken.	Pas un proverbe
This is what happens when we start out on the wrong path.	Pas un proverbe
This is my favourite drink.	Pas un proverbe

Tableau 3. – En partant de « This »

Intéressants 0

Etrange 0

Pas un proverbe 10

Incohérents 0

Le modèle de base étant plus important et donc déjà plus entraîné, il ne crée pas de résultat incohérent comme le précédent, mais le dataset utilisé est trop petit pour créer des proverbes intéressants. Peu de proverbes commencent par « Some » dans le dataset. ($\frac{4}{3208} \approx 0.12\%$)

4.1.2. Entraînement sur les proverbes traduits

4.1.2.1. Modèle de départ facebook/opt-125m

On entraîne d'abord sur **20000** proverbes (57.14%).

Proverbe	Commentaire
Some men are good at it, some are bad at it, and some are good at it.	Incohérent
Some day people will get one.	Pas un proverbe
Some days a man wears a leather jacket and a woman wears a leather garment.	Pas un proverbe
Somehow the worst person in the world can make the best at things.	Étrange
Somehow, I am not a thief.	Pas un proverbe
Some people are more beautiful than others.	Pas un proverbe
Some day, the sun will rise on a mountain.	Incohérent
Someones got a good idea and they're good at it.	Pas un proverbe
Some men have a heart to die for.	Intéressant
Somehow it looks like a human's ear.	Incohérent

Tableau 4. – En partant de « Some »

Proverbe	Commentaire
A dog is a beast of a mind.	Incohérent
A man who makes his wife cry is a thief.	Étrange
A man's heart is not a dog's tongue.	Incohérent
A girl may be a good girl, but she has a better chance of getting married.	Étrange
A man who has never seen a woman is a man.	Incohérent
A house that is just a house is not a house that is a house.	Incohérent
A man who will not sleep for an hour, will not sleep for an hour.	Incohérent
A man cannot take a wife.	Pas un proverbe
A man who eats meat will not be a judge.	Étrange
A bit of good luck in life is better than nothing.	Intéressant

Tableau 5. – En partant de « A »

Intéressants 2

Etrange 4

Pas un proverbe 6

Incohérents 8

Sur cette nouvelle génération à partir de 20000 proverbes une majorité de proverbes non intéressants ont été générés. Sur d'autres générations (non représentées ici) les résultats varient avec plus de proverbes intéressants. Cela peut être dû au mot de départ de la génération qui peut avoir trop de possibilités ou pas assez et qui essaie de déterminer le prochain mais sans grande cohérence.

4.1.2.2. Modèle de départ TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T

La RAM disponible ne nous a permis que de sélectionner un maximum de 5000 proverbes.

Cependant, on obtient quand même des résultats plus intéressants qu’avec 3208 proverbes

Proverbe	Commentaire
Some men are born to lead, and some to follow.	Intéressant
Some men are so proud of their looks that they never look at the rest of their faces.	Étrange
Someone will not be able to find a wife whom his father chooses.	Pas un proverbe
Some things are good for the body, but not for the stomach.	Intéressant
Some are born with the gift of knowledge, and some with the gift of ignorance.	Intéressant
Some folks make their own pies, and some do not.	Pas un proverbe
Some people call the wind the sun’s enemy.	Incohérent
Some are wiser than they know. Humor as a way of life is also a part of the world I have to live in.	Incohérent
Some are born great, but some become great by their wit.	Intéressant
Someone is not what you think about him; you think about him.	Incohérent

Tableau 6. – En partant de « Some »

Proverbe	Commentaire
A man’s name is his life.	Intéressant
A man who has never taken a step should never be trusted.	Intéressant
A great king cannot be a great man.	Intéressant
A woman is hard to handle, but easy to cheat.	Étrange
A dog’s nose is better than a man’s eyes.	Intéressant
A woman who has a good head cannot be a fool.	Intéressant
A bird from a city is a nesting place for many.	Incohérent
A little knowledge is better than a great ignorance.	Intéressant
A man who speaks of a hundred will be thought a hundred.	Incohérent
A man should not be too ambitious.	Intéressant

Tableau 7. – En partant de « A »

Intéressants 11

Etranges 2

Incohérents 5

Pas un proverbe 2

Comme précédemment, ce modèle produit peu de résultat incohérents. On constate qu’on obtient des résultats plus intéressants avec ce dataset plus important, plus adapté à la taille du modèle de base.

4.1.3. Expérience avec des humains

A titre d'expérience, nous avons rassemblé 24 proverbes originaux et générés par nos modèles. (facebook/opt-125m sur ≈ 15000 proverbes) Nous avons demandé à 18 personnes (des professeurs de SHS notamment) de classer chaque proverbe à l'aveugle.

Proverbe	Généré	Taux Succès
The horse that is not ridden is worse than the horse that is not ridden.	Oui	100%
Set a fire and burn it.	Oui	61.11%
Think and thank god.	Non	44.44%
The best way to see divine light is to put out your own candle.	Non	66.67%
The riches of the mind may make a man rich and happy.	Non	50%
Who pays in advance eats stinking fish.	Non	33.33%
The nail that sticks out gets pounded.	Non	66.67%
Work paid for in advance has feet of lead.	Non	22.22%
A young child is not born in spite of his own mother.	Oui	61.11%
Set your own goals. An arrow, a sword, and a bottle of wine are enough to keep you from leaving your wife.	Oui	72.22%
All men can't be first.	Non	77.78%
Don't celebrate until you are across the brook.	Non	55.56%
Such carpenters, such chips.	Non	11.11%
A person's only hope is to give him a good laugh.	Oui	83.33%
To pay one back in one's own coin.	Non	33.33%
An English friend is a friend.	Oui	83.33%
Romeo must die in order to save the love.	Non	16.67%
The only thing that can be done is to make them believe that they are not going to die.	Oui	50%
To be or not to be, you have to be.	Oui	77.78%
Children are poor men's riches.	Non	16.67%
He, that would eat the fruit, must climb the tree.	Non	88.89%
Do not drink in advance on the hide of a bear.	Non	16.67%
A man's clothing can only last one lifetime.	Oui	50%
The musician who is paid in advance does not play so well.	Non	50%

Tableau 8. – Résultats de l'expérience par proverbe

	Générés	Non générés	Total (%)
Bien détectés	115	117	53.704
Mal détectés	153	47	46.296
Total (%)	37.5	62.5	100

Tableau 9. – Résultats de l'expérience

4.2. Génération avec thème

On a ici entraîné le modèle facebook/opt-125m sur 20000 proverbes (57.14%) annotés avec leurs thèmes (1 heure d'entraînement).

Proverbe	Commentaire
The love of one is not always the love of all.	Intéressant
If you love a woman, she will be love.	Etrange
If you love a man, you will be love.	Etrange
Love is love, and with love we can go out of love and into another.	Etrange
Love is the best proverb.	Intéressant
It is a good proverb, but a proverb is better.	Incohérent
Love has no love.	Etrange
It is love that first love is in the end.	Intéressant
Never in love. It's not love.	Intéressant
Love the love of the love of peace, and the love of the love of the love of the love of the peace...	Incohérent
The woman heches the dog, but the man sheches the dog.	Etrange
If love can't be found or found without fear, it can't be found without fear.	Etrange

Tableau 10. – Avec le thème « Love »

Intéressants 4

Etrange 6

Pas un proverbe 0

Incohérents 2

Le plus long temps d'entraînement et la rajout de contexte a permis d'obtenir uniquement des proverbes malgré l'étrangeté ou l'incohérence de certains. On remarque que lors d'ajout d'un thème, la génération essaie souvent d'inclure le mot dans le proverbe généré.

5. Conclusion

Dans l'ensemble, notre modèle renvoie des proverbes satisfaisants à la génération (que ce soit avec facebook ou TinyLlama). Dans certains cas les proverbes sont incohérents ou étranges. Cela provient de *notre interprétation d'un proverbe* car en regardant des proverbes de notre bdd, certaines générations sont ressemblantes, provenant probablement des traductions vers l'anglais pouvant perdre en précision.

La difficulté principale rencontrée est la **validation du modèle** car vérifier que les générations sont « correctes » est assez subjectif. Mais nous avons quand même pu différencier ceux « mauvais » de ceux de « meilleurs » qualité en tant que proverbes.

Réaliser un sondage auprès de professeurs de SHS a aussi été un bon moyen de déterminer la **précision** du générateur car dans presque un **tiers des cas** il y avait une confusion entre vrai et généré.

Pistes d'améliorations :

- Trouver plus de proverbes issus d'un anglais natif et non traduit → meilleure précision dans la génération possible (avoir moins de vérité générale ou autre),

- Tester différents réseaux et déterminer lequel serait le plus adapté à notre problème de génération.