

Rapport Amal 3

Mathieu GROSSO - Thomas Floquet

December 2021

1 TP07 - Régularisation et protocole expérimental

1.1 Préparation du modèle:

On entraîne un réseau de neurones simples avec 3 couches linéaires de 100 channels. La sortie est de taille 10 puisqu'on a 10 classe (chiffres de 0 à 9). On utilise une cross entropy pour entraîner le modèle, on optime avec Adam sans lr decay. On choisit un Lr de 0.001, 1000 epochs, et des batchs de 300 données. De même on ne conserve que 0.05 des données de train et de val.

Sans aucune régularisation, on voit que le modèle apprend vite en train et en test mais que rapidement il overfit et que la courbe de train en test augmente après 100 epochs jusqu'à remonter à son niveau de départ. De même l'accuracy diminue. On va visualiser la loss, les poids et le gradient pour essayer de voir l'évolution du modèle au cours de l'entraînement.

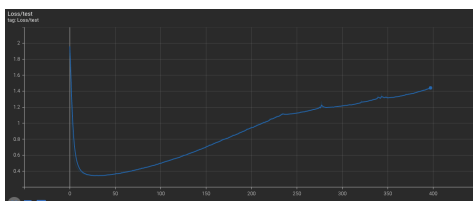


Figure 1: Loss en test sans aucune régularisation, avec un learning rate de 0.0001.

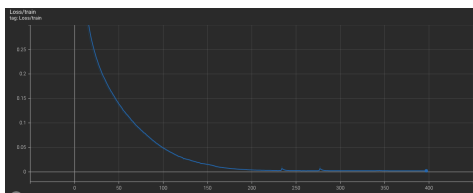


Figure 2: Loss en train sans aucune régularisation, avec un learning rate de 0.0001.

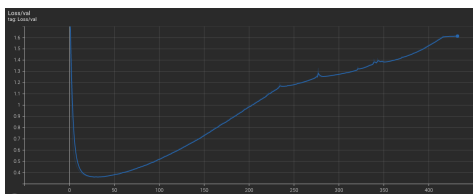


Figure 3: Loss en validation sans aucune régularisation, avec un learning rate de 0.0001.

En train l'apprentissage se passe bien, on apprend bien. Le problème est que sans régularisation on finit par sur apprendre les données de train et que la loss remonte en test et en validation. On se retrouve donc avec de mauvais résultat en test et en validation comme on peut le voir sur le graph.

La loss qui descend jsuqu'à 0.2 après 50 epochs remonte à 2 en validation.

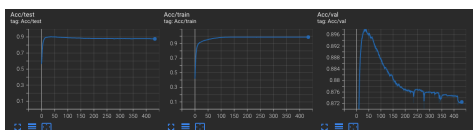


Figure 4: Accuracy en test, val et train sans aucune régularisation.

Le graphe d'accuracy montre la même chose que la courbe de loss, on a surtout une très mauvaise accuracy en validation après 100 epochs. On rajoute également les histogrammes de visualisation des poids et des valeurs du gradient en training, test et validation.

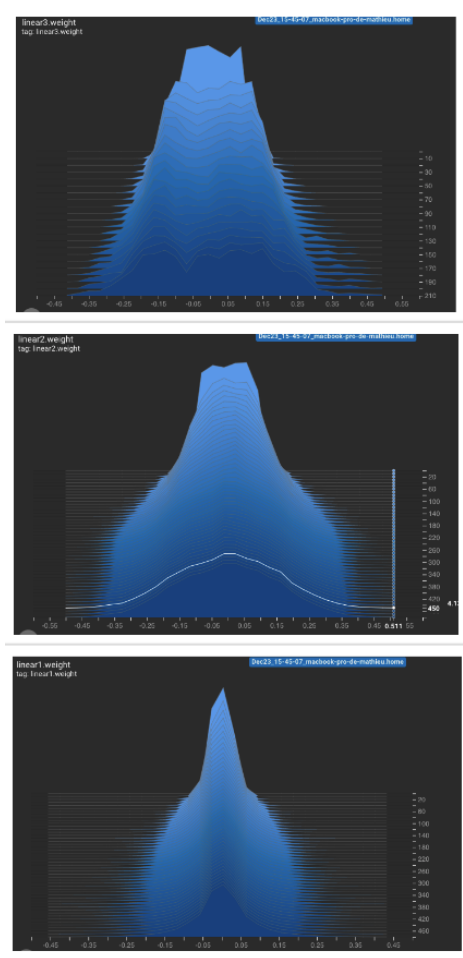


Figure 5: Différents poids de chaque couche au cours de l'apprentissage sans régularisation.

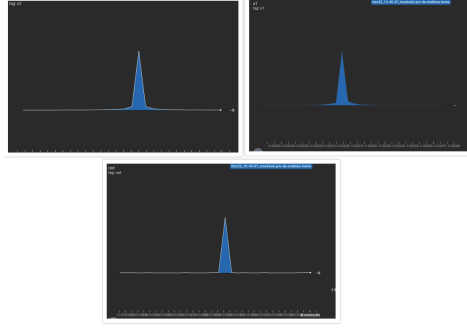


Figure 6: Gradient au cours de l'apprentissage sans régularisation.

1.2 Régularisation:

On va ensuite s'intéresser à la régularisation du modèle; On va tester différentes approches, le dropout, la régularisation L1 et L2, la batchnorm et la layer norm. Au final, on retiendra le dropout, la régularisation L1 et un decay du learning rate de 0.009. En rajoutant cette régularisation, on a des résultats qui convergent plus vite et surtout le modèle ne sur-apprend pas en train puisque même après 400 epochs la loss en val et en test en remonte pas.

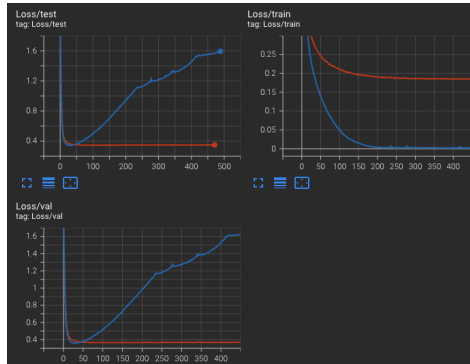


Figure 7: Loss en train, test et validation avec régularisation (dropout, L1, lr decay).

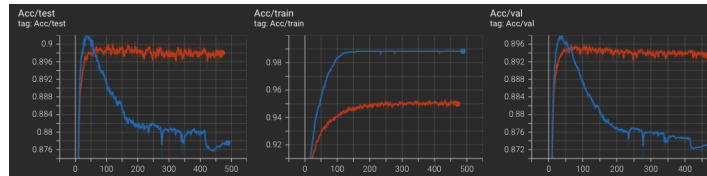


Figure 8: Accuracy en train, test et validation avec régularisation (dropout, L1, lr decay).

En rouge la courbe avec régularisation ($dropout = 0.47$, $lambda = 0.001$ (pour la régularisation L1), et $decay = 0.009$). Sur la courbe de la loss, on observe clairement une amélioration en val et en test. Contrairement à avant où le modèle ne convergeait pas, maintenant il converge et plus haut que sans la régularisation. C'est particulièrement vrai en validation où la différence est importante. De même sur les courbes d'accuracy on a une bien meilleure accuracy en val et en test. A l'inverse en train, les scores sont légèrement moins bons mais c'est normal puisqu'on sur apprend moins sur les données de train.

Ensuite on va comparer les poids et gradients avec régularisation et sans régularisation.

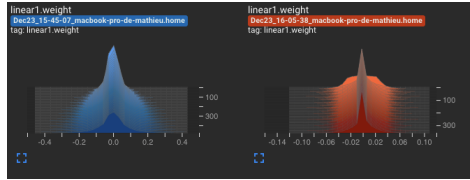


Figure 9: Différents poids de chaque couche au cours de l'apprentissage avec régularisation.

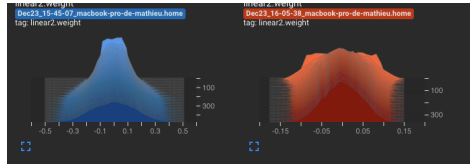


Figure 10: Différents poids de chaque couche au cours de l'apprentissage avec régularisation.

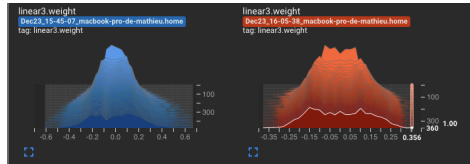


Figure 11: Différents poids de chaque couche au cours de l'apprentissage avec régularisation.

En rouge les poids pour l'apprentissage avec régularisation. On voit que le modèle en rouge semble converger mieux.

1.3 Augmentation de données:

l'augmentation des données réduit grandement l'accuracy. Les opérations effectuées sont: $data_transforms2 = transforms.Compose([transforms.ToPILImage(), transforms.RandomRotation(degrees = 45), transforms.Random(18, 18), transforms.ToTensor()])$.

On a donc décidé de ne pas garder cela dans la suite.

1.4 Optuna, Lightning, et optimisation des hyperparamètres:

On a effectué une grid search de 100 itérations sur le coefficient de régularisation L1 et sur le coefficient de dropout. Pour le dropout on a placé le coef entre 0.1 et 0.5 et pour la régularisation L1 on a choisi le coefficient entre 0.0001 et 0.1. Puis on a optimisé grace à optuna en créant un objectif et une étude sur le modèle précédemment utilisé.

Les résultats sont les suivants:

2 TP08 - Détection de sentiments :

Nous avons étudié 3 architectures différentes :

1. 1 couche de convolution avec 16 out_channels
2. 2 couches de convolution allant de 16 à 32 out_channels

3. 4 couches de convolution allant de 16 à 128 out_channels

Tous suivis d'un MLP à deux couches.

Nous analysons la performance relative de notre modèle en comparant l'accuracy à un modèle renvoyant toujours la classe majoritaire 0.

Question 3 :

$$S(i+1) = s(i+1) * s(i) \quad (1)$$

$$W(i+1) = w(i+1) + k(i-1) \quad (2)$$

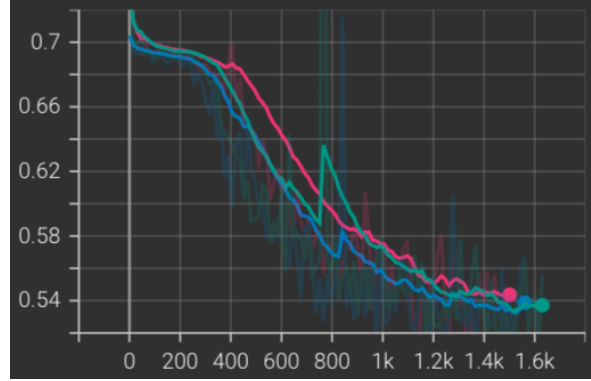


Figure 12: Loss of the 3 models

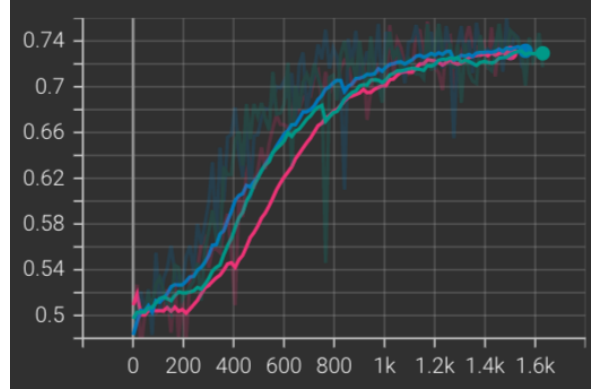


Figure 13: Accuracy of the 3 models

Légende : En vert le modèle 1, en bleu le modèle 2 et en rose le modèle 3.

On remarque que les 3 modèles ont des performances assez proches, avec un léger avantage pour le modèle 2 qui semble converger un peu plus vite et avoir une très légère meilleure performance. L'accuracy relative montre bien que nos modèles sont plus performants qu'un modèle renvoyant toujours la classe majoritaire (ici ayant une accuracy d'environ 0.5).

3 TP09 - CNN et classification de sentiments

BasicNet : modèle simple de la question 1 où chaque mot est utilisé avec la même pondération.

SimpleAttention : modèle d'attention de la question 2 qui pondère les mots qui ont un intérêt pour la classification. Une fonction simple qui permet d'obtenir $p(a_i|t)$ est la fonction softmax

ComplexAttention : modèle de la question 3 où la query dépend du texte lui-même.

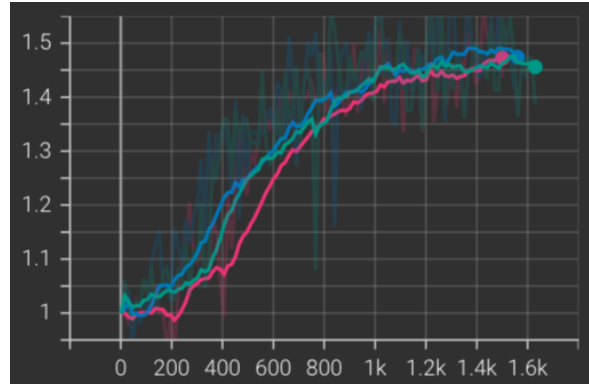


Figure 14: Relative accuracy of the 3 models of TP8

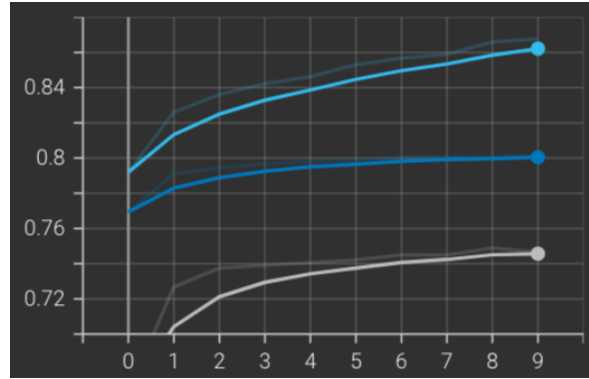


Figure 15: Accuracy of the 3 models of TP9

Légende : En gris le modèle BasicNet, en bleu foncé le modèle SimpleAttention et en bleu clair le modèle ComplexAttention.

Comme prévu, ComplexAttention obtient de meilleurs résultats que SimpleAttention qui lui-même obtient de meilleurs résultats que BasicNet (Figure 16). Cela s'explique par le fait que le modèle ComplexAttention apprend la query de manière à la faire dépendre du texte lui-même, ce qui permet d'avoir une query plus intéressante pour l'analyse de sentiment.

On remarque que le modèle s'intéresse bien aux mots les plus importants pour l'analyse de sentiments. En effet, les mots avec les valeurs d'attention les plus élevées sont *hopelessly* et *goeey*, et les mots avec les valeurs les moins importantes sont *hard* et *correctness*.

4 TP10 - Attention globale

4.1 1 - Modèle de base

Nous avons d'abord étudié un modèle simple de transformers avec $L = 3$. (3 modules de self Attention). Nous avons utilisé 100 pour la dimension des embeddings, un learning rate de 0.001, 10 époques et une batch size de 128. Le modèle apprend vite et converge autour de 0.75 ce qui est meilleur que ce que nous avons trouvé dans les TP précédents (cf TP09 au dessus). Les hyperparamètres pour les dimensions sont les suivants et sont disponible dans le fichier config du code:

Ensuite nous avons testé en rajoutant toute les options d'un transformers

```

input phrase: ['somewhere', 'in', 'the', 'dark', 'recesses', 'of', 'my', 'brain', 'cells', 'a', 'song', 'plays', 'in', 'my', '___00V___',
'i', '___00V___', 'forget', 'it', 'no', 'matter', 'how', 'hard', 'i', '___00V___', '___00V___', 'midnight', 'madness', 'and', '___00V___',
'gonna', 'get', 'to', '___00V___', 'wish', 'i', 'could', 'find', 'a', 'copy', 'of', 'this', 'on', 'a', '45rpm', '___00V___', 'five',
'disparate', 'teams', 'head', 'out', 'one', 'night', 'in', 'l.a.', 'for', 'a', 'scavenger', 'hunt', 'for', 'clues', 'instead', 'of',
'physical', '___00V___', 'an', 'unkempt', '___00V___', 'with', 'two', 'gorgeous', 'assistants', 'is', 'the', 'mastermind', 'of', 'all',
'this', 'insanity', '___00V___', 'about', 'to', 'be', 'unleashed', 'on', 'l.a.', 'all', 'the', 'teams', 'are', 'stereotypes', '___00V___',
'movie', 'being', 'from', '___00V___', 'before', 'political', 'correctness', 'screwed', 'everything', '___00V___', 'the', '___00V___',
'___00V___', 'the', '___00V___', 'led', 'by', 'eddie', '___00V___', 'the', 'dumb', '___00V___', '___00V___', 'the', '___00V___', '___00V___',
'especially', 'the', '___00V___', 'the', 'giggling', 'twins', 'are', 'a', '___00V___', '___00V___', 'and', '___00V___', 'the', '___00V___',
'___00V___', 'with', 'stephen', 'furst', 'as', 'the', '___00V___', 'furst', 'is', 'hilarious', 'as', 'the', 'overweight', 'slob', '___00V___',
'whose', 'attempt', 'to', 'use', 'a', 'computer', 'to', 'decipher', 'the', 'various', 'clues', 'leads', 'to', 'a', 'goosey', '___00V___',
'movies', 'like', 'this', '___00V___', 'made', '___00V___', 'these', '___00V___', 'movies', 'have', 'to', 'have', 'an', '___00V___', 'to',
'them', 'with', 'some', 'dark', 'characters', 'and', 'other', '___00V___', 'go', 'back', 'to', 'the', 'days', 'when', 'the', '___00V___',
'___00V___', 'led', 'by', 'david', 'naughton', 'were', 'still', 'good', 'and', 'not', 'hopelessly', '___00V___', 'so', 'dump', 'all',
'serious', 'pretensions', 'and', 'go', 'back', 'to', '___00V___', '___00V___', 'midnight', 'madness', '.', '.', '.']
word: hopelessly, attention value: 0.040796324610710144
word: goosey, attention value: 0.03663202375173569
word: slob, attention value: 0.022535987198352814
word: serious, attention value: 0.020592158660292625
word: hilarious, attention value: 0.019024258479475975
word: unleashed, attention value: 0.016735348850488663
word: good, attention value: 0.016594644635915756
word: gorgeous, attention value: 0.015835488215088844
word: especially, attention value: 0.01349448412656784
word: pretensions, attention value: 0.012700004503130913
word: dumb, attention value: 0.012213854119181633
word: forget, attention value: 0.012088504619896412
word: unkempt, attention value: 0.011397838592529297
word: hard, attention value: 0.010123910382390022
word: correctness, attention value: 0.010118484497070312

```

Figure 16: Attention values du TP9

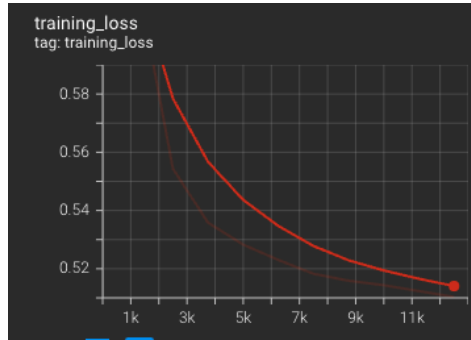


Figure 17: loss en train TP10

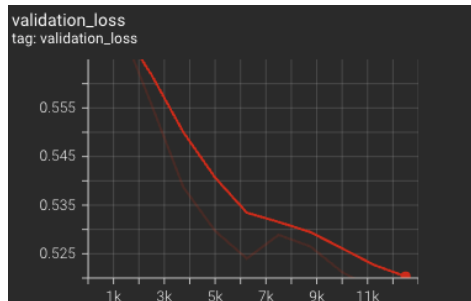


Figure 18: loss en val TP10

4.2 2 - Ajout d'un modèle Résiduel, d'un token de CLS et d'un positional encoding

Dans cette section, nous avons rajouté des connexions résiduel, un token CLS et le positional encoding. Les hyperparamètres des dimensions sont les mêmes que précédemment ainsi que les paramètres d'entraînement. En revanche nous avons fait plus d'époques car au bout de 10 époques les résultats étaient moins bon que précédemment. En effet l'ajout d'un modèle pour le token peut ralentir le

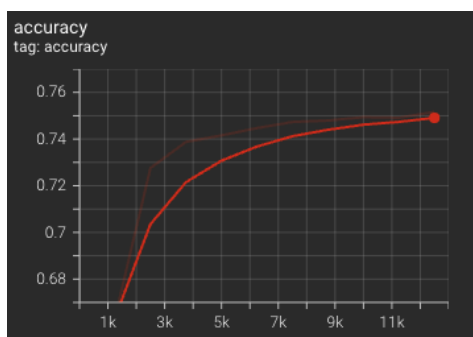


Figure 19: Accuracy en train.

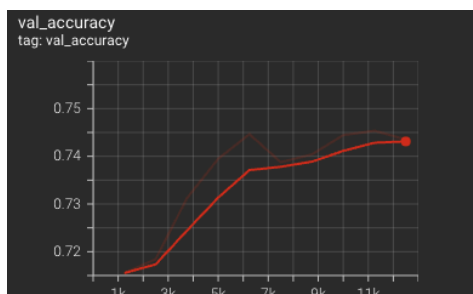


Figure 20: Accuracy en val.

training.

Au bout de 30 epoques voilà les résultats: sur les graphiques en rouge ce sont les résultats pour le modèle basique de transformers sans les ajouts de connections résiduelles, de token et d'encoding, et en vert avec:



Figure 21: loss en train.

On observe que la courbe rouge obtient de meilleurs résultats et semble aussi plus stable. Nous avons aussi testé avec plus d'époques pour l'encodeur basique et en effet il semble mieux fonctionner. Nous avons aussi essayé de retirer le token cls et on obtient des résultats comparables. On ne met pas les courbes car ils sont très proches.

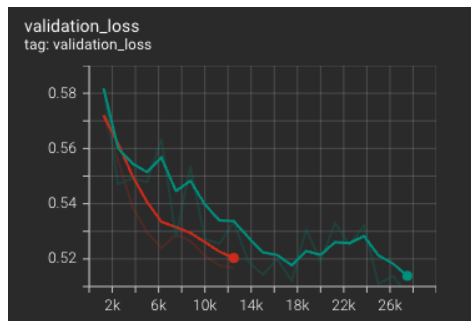


Figure 22: loss en val.

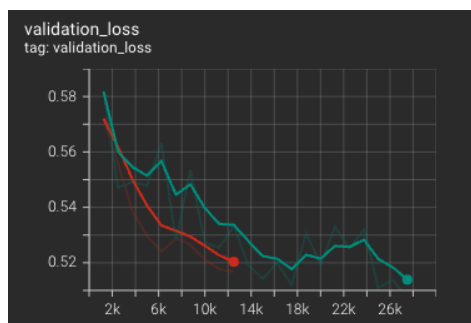


Figure 23: Accuracy en train.

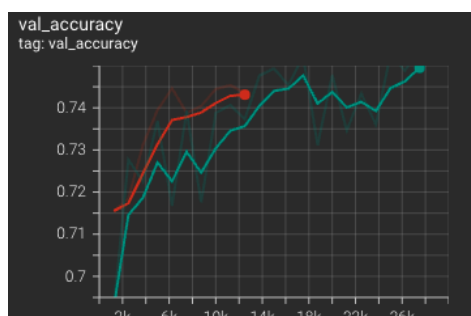


Figure 24: Accuracy en val.