

Rapport TME1 : Regression

CELLIER Roxane - MOULIN-ROUSSEL Lou

1

1. Batch Linear Least Squares

Pour étudier les différences entre les deux premières méthodes proposées, nous avons effectué 2000 tests sur différents jeux de données de même taille (50), pour nous permettre d'évaluer les temps d'exécution et les écarts entre les paramètres appris. La boucle que nous avons utilisée est la suivante :

```
1  for i in range(2000):
2      batch = Batch()
3      size = 50
4      batch.make_linear_batch_data(size)
5      model = LinearModel(size)
6      x = np.array(batch.x_data)
7      y = np.array(batch.y_data)
8
9      start = time.process_time()
10     model.train(x, y)
11     theta1[i, :] = model.theta
12     time_lls[i] = time.process_time() - start
13
14     start = time.process_time()
15     model.train_from_stats(x, y)
16     theta2[i, :] = model.theta
17     time_lls_stats[i] = time.process_time() - start
18
```

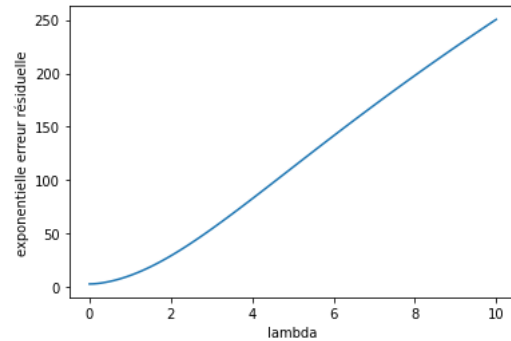
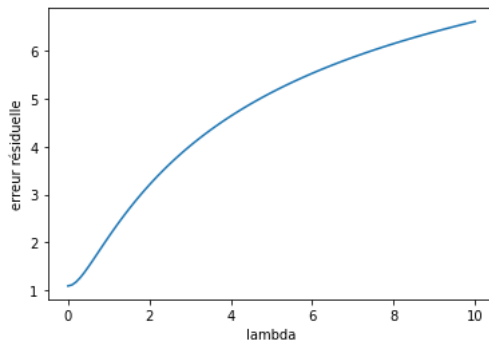
Ce qui nous a permis d'obtenir les résultats suivants en calculant les moyennes :

```
différences de theta : [6.66133815e-16 0.00000000e+00]
moyenne de temps lls : 0.00012791794249999988
moyenne de temps lls avec stats : 0.00046454361200001413
```

On semble constater que l'utilisation de stats est légèrement moins rapide, même si les deux méthodes s'exécutent en un temps du même ordre grandeur. De plus, les valeurs renvoyées des paramètres theta ne sont pas forcément exactement identiques, mais cependant extrêmement proches. Cela peut être dû à des approximations et arrondis faits par la machine lors des calculs.

2. Ridge Regression

Pour ce test, nous avons considéré un jeu unique de taille 50, pour lequel nous avons testé différentes valeurs du paramètre lambda, de 0 à 10 par pas de 0,1, pour ainsi observer la variation de l'erreur.



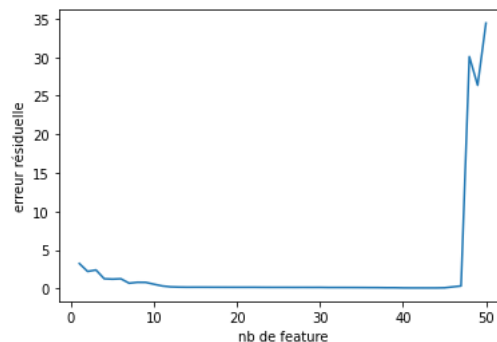
Nous remarquons que l'erreur résiduelle croît lorsque le poids λ augmente. Elle semble augmenter de manière logarithmique. En effet, sur la deuxième figure ci-dessus, nous pouvons voir que l'exponentielle de l'erreur suit une progression affine voire linéaire en fonction de λ .

3. Radial Basis Function Networks

De la même manière que dans la première partie (Batch Linear Least Squares), nous avons effectué 2000 tests sur des jeux différents pour pouvoir comparer en moyenne le temps d'exécution des deux approches présentées dans ce TME. Nous constatons assez vite que les deux perspectives ont des temps d'exécution très similaires, la première approche semblant tout de même sensiblement plus rapide que la deuxième.

moyenne de temps RBNF1 : 0.0005012852634999945
moyenne de temps RBNF2 : 0.0007240960204000352

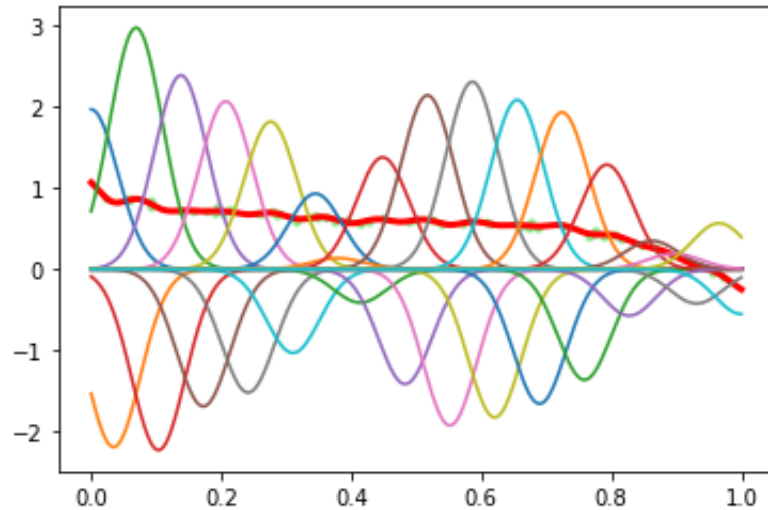
Nous regardons maintenant l'évolution de l'erreur sur un même jeu de données de taille 60, selon le paramètre λ . Nous avons choisi ici d'utiliser la première perspective, celle-ci étant sensiblement plus rapide.



L'erreur résiduelle décroît progressivement plus on augmente le nombre de features. On observe une diminution très forte entre 0 et 5 features et une diminution qui ne semble plus significative au-dessus de 10 features. Cependant sur certains tests, l'écart remonte fortement autour de 50 features, comme illustré par la figure ci-dessus.

Avec 30 features, on obtient le graphe ci-dessous, avec en rouge et en gras la courbe

obtenue par régression et en vert les points avec lesquels la régression est faite. Les autres courbes représentent les features.



Nous constatons ici un overfitting du modèle sur les données de test : la courbe formée est très proche de chacun des points et nécessite de nombreux paramètres pour être décrite.

4. Incremental RBFNs

Nous allons ici tenter de comparer globalement chacune des méthodes pour estimer ce qui rend une approche meilleure qu'une autre. Nous avons voulu dans un premier temps constater les éventuelles variations liées à la taille du jeu de données. Nous avons donc pris différentes tailles allant de 50 à 500, sur lesquelles nous avons effectué nos tests sur 1000 jeux différents pour nous permettre de fournir une moyenne. Nous avons choisi pour ces tests de fixer le nombre de feature à 15 et le nombre maximal d'itération à 500. Voici le code effectuant les 1000 tests pour une taille t spécifique :

```

1  for i in range(NB_TRY):
2      # cr ation des jeux de donn es
3      batch = Batch()
4      batch.make_nonlinear_batch_data(size=t)
5      x = np.array(batch.x_data)
6      y = np.array(batch.y_data)
7
8      # test temps GDRBFN
9      modelGD = GDRBFN(nb_features=NB_FEATURE)
10     start = time.process_time()
11     for j in range(MAX_ITER):
12         x, y = batch.add_non_linear_sample()
13         modelGD.train(x, y, alpha=0.5)
14     time_GDRBNF[i] = time.process_time() - start
15     error_GDRBNF[i] = modelGD.compute_error(batch.x_data, batch.y_data)
16
17     # test temps RLSRBFN
18     modelRLS = RLSRBFN(nb_features=NB_FEATURE)
19     start = time.process_time()

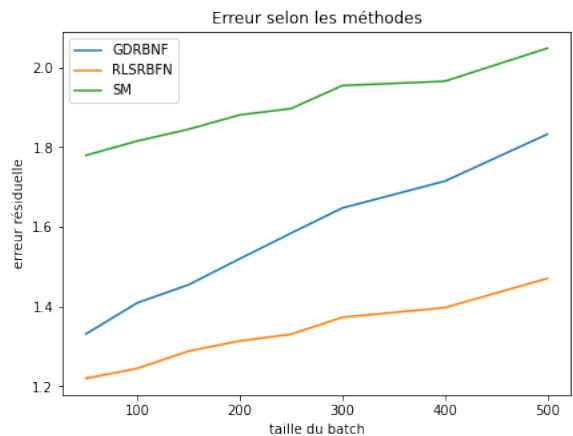
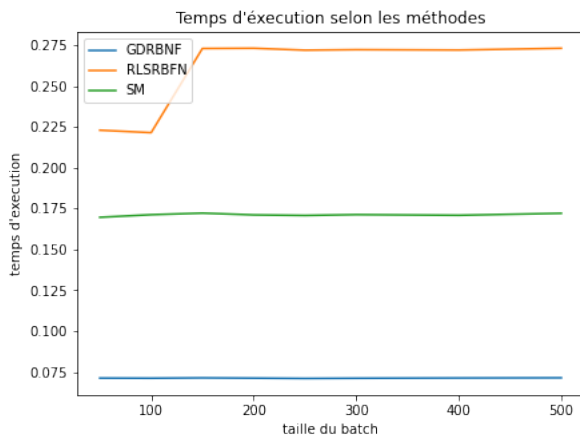
```

```

20     for j in range(MAX_ITER):
21         x, y = batch.add_non_linear_sample()
22         modelRLS.train(x, y)
23         time_RLSRBFN[i] = time.process_time() - start
24         error_RLSRBFN[i] = modelRLS.compute_error(batch.x_data, batch.y_data)
25
26     # test temps Sherman - Morrison
27     modelSM = RLSRBFN(nb_features=NB_FEATURE)
28     start = time.process_time()
29     for j in range(MAX_ITER):
30         x, y = batch.add_non_linear_sample()
31         modelSM.train_ribs_sherman_morrison(x, y)
32         time_SM[i] = time.process_time() - start
33         error_SM[i] = modelSM.compute_error(batch.x_data, batch.y_data)
34

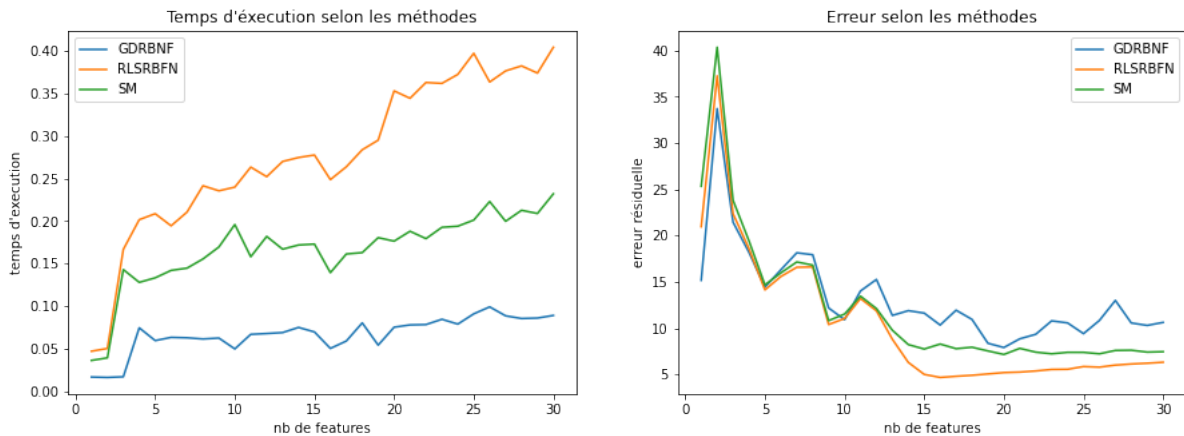
```

Ce qui nous donne comme affichage :



Les temps d'exécution selon les méthodes semblent constants en fonction de la taille du batch, excepté pour la méthode Recursive Least Squares qui semble avoir un léger sursaut autour d'un jeu de 100 valeurs. Les erreurs selon les méthodes semblent quant à elles augmenter légèrement en fonction de la taille du batch. Cela peut paraître logique puisque la taille des données augmente, mais pas le nombre de features permettant de coller à ces données.

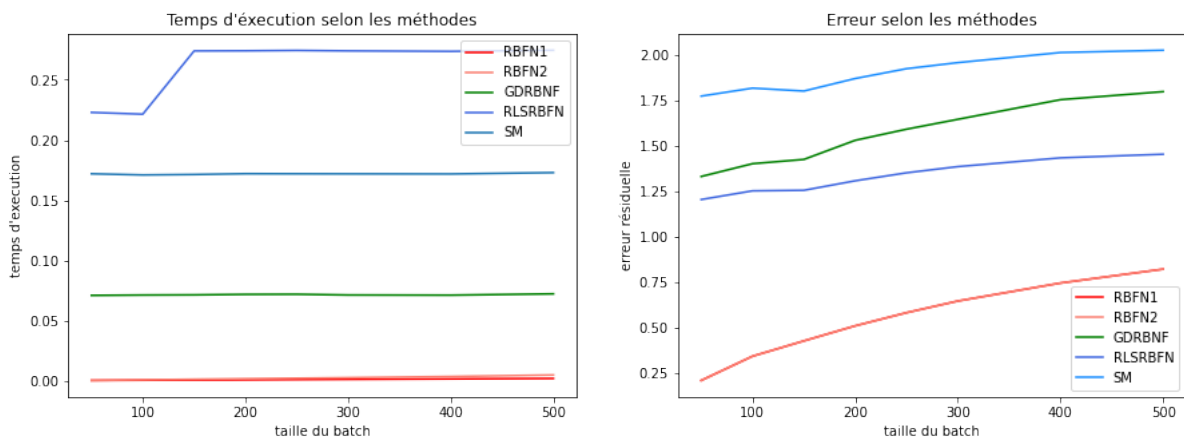
Nous cherchons maintenant à faire varier le nombre de features pour un même jeu de données, puis nous affichons les temps d'exécution et les erreurs selon les méthodes.



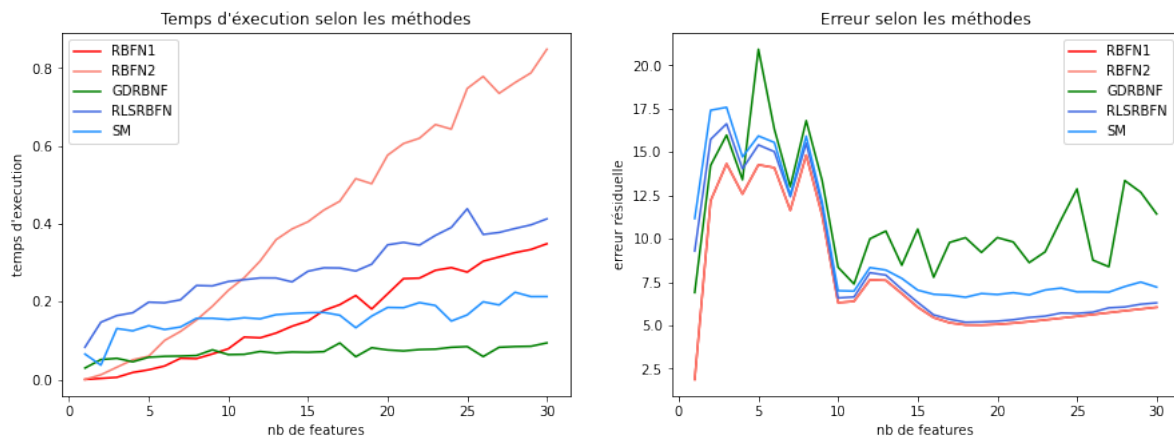
On observe que la méthode de descente de gradient semble plus rapide mais l'erreur est plus grande lorsque le nombre de features est supérieur à 10 environ. La méthode de recursive least squares sans l'approche de Sherman-Morrison (RLSRBFN) est la plus lente mais aussi celle avec la plus petite erreur.

5. Comparaison Batch et Incremental

En utilisant les mêmes types de test et les mêmes paramètres (nombre de tests, tailles des jeux, etc) sur l'ensemble des méthodes RBNF, nous obtenons les mêmes types de tracé. Pour simplifier la visualisation, les méthodes Batch sont teintées en rouge et les méthodes Incremental sont teintées en bleu.



Nous pouvons voir de manière assez nette que les méthodes batch semblent nettement plus rapide que les méthodes incrémentales, elles semblent également renvoyer moins d'erreur que les autres méthodes.



On peut remarquer que pour un faible nombre de features, les techniques incrémentales sont plus lentes et ont une plus large erreur résiduelle. Cependant si l'on augmente grandement le nombre de features, on peut voir que le temps d'exécution croît beaucoup plus vite et en vient même à dépasser celui des méthodes incrémentales. De manière globale, la méthode du gradient descendant semble la plus rapide, en effet son temps d'exécution ne varie presque pas en fonction du nombre de features, mais comme précédemment, son erreur est la plus grande pour un large nombre de features.

Au vue des courbes visibles, il n'est pas évident de décider qu'une méthode est largement supérieure à toutes les autres. Cela va fortement dépendre de ce que l'on cherche à minimiser (temps ou erreur), de la taille des données, et également du nombre de features que l'on souhaite récupérer pour notre modèle.

6. Locally Weighted Regression

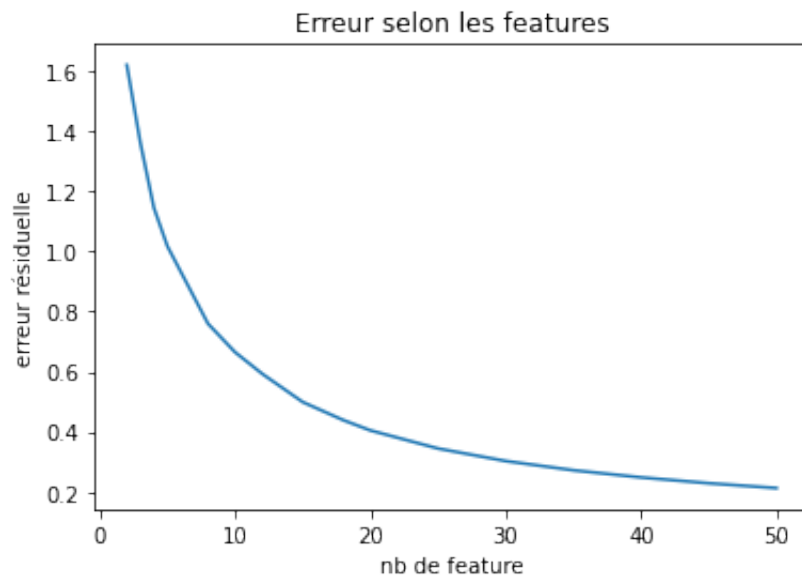
Nous allons ici visualiser l'impact du nombre de sous-divisions sur l'erreur du modèle appris par LWR. Pour cela, nous allons prendre un jeu de données fixe sur lequel nous allons apprendre différents modèles selon le nombre de sous-ensembles locaux voulus. Nous testerons l'erreur pour différents découpages de 2 à 50 ensembles. Le code nous permettant d'effectuer ces tests est présenté ci-dessous :

```

1  batch = Batch()
2  batch.make_nonlinear_batch_data()
3
4  nb_feature = [2,3,4,5,8,10,12,15,18,20,25,30,35,40,45,50]
5  error = np.zeros(len(nb_feature))
6
7  for i in range(len(nb_feature)):
8      model = LWR(nb_features=nb_feature[i])
9      model.train(batch.x_data, batch.y_data)
10
11     error[i] = model.compute_error(batch.x_data, batch.y_data)
12

```

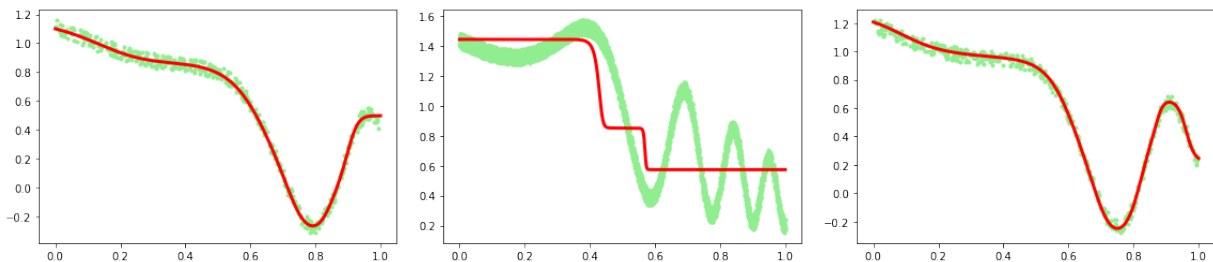
Suite à l'exécution du code et à l'affichage, nous obtenons le tracé suivant :



De façon assez logique, on peut voir que l'erreur baisse au fur et à mesure que l'on divise le jeu en plusieurs sous-ensembles. Il semble en effet plus efficace pour essayer de suivre une courbe de faire plusieurs petits segments plutôt qu'une seule droite.

7. Regression with Neural Networks

Nous allons observer différentes utilisations de PyTorch pour apprendre différents modèles autour de différents jeux, et étudier les différents résultats ainsi que le respect des données par le modèle ainsi appris.



On peut constater que le modèle résultant de la méthode Batch (à droite) et celui résultant de la méthode incrémentale utilisant des mini-batches (à gauche) sont cohérents avec le jeu de données. Cependant la version incrémentale (au milieu) ne donne pas un résultat concluant. Le modèle proposé ressemble plus à une fonction en escalier, et les marches semblent représenter une moyenne générale sur la largeur qu'elles représentent.

Nous allons ensuite étudier l'évolution du temps d'exécution et de la valeur de la loss suivant la dernière méthode, en faisant varier les tailles des mini-batches. Nous avons utilisé cette boucle nous permettant d'effectuer nos tests sur une taille s :

```

1  for t in range(nb_try):
2      start = time.process_time()
3      for i in range(max_iter):
4          xt, ym = make_minibatch(minibatch_size, batch.x_data, batch.y_data,
                                  size)

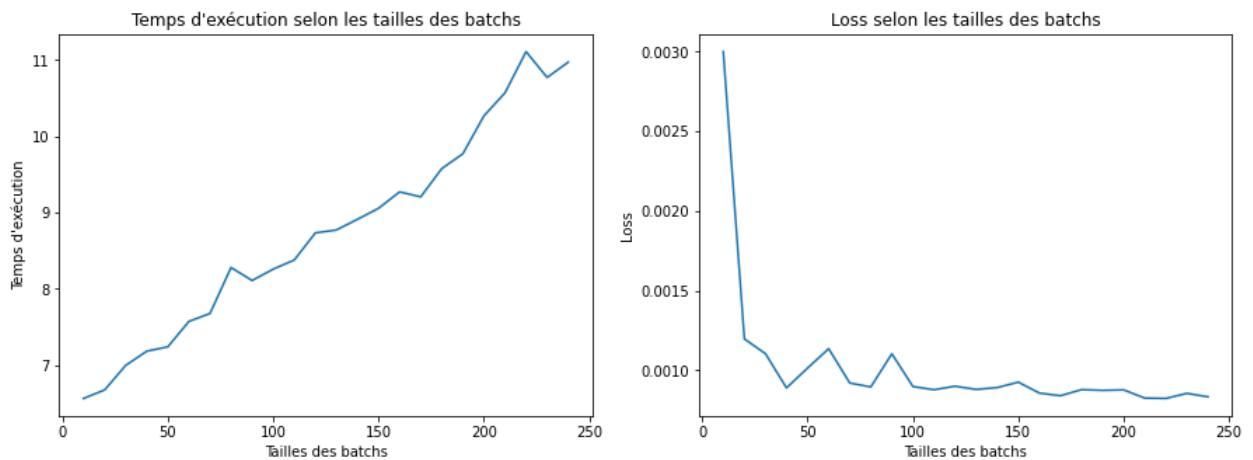
```

```

5         yt = th.from_numpy(ym).float()
6         output = model.f(xt)
7         loss = func.mse_loss(output, yt)
8         #print(loss.item(), model.compute_error(xt,yt))
9         model.update(loss)
10        #print("NN Incr Reg time:", time.process_time() - start)
11        times_s[t] = time.process_time() - start
12        loss_s[t] = loss
13

```

Nous avons utilisé la valeur `nb_try = 10` pour avoir une moyenne sur 10 valeurs pour chaque taille de minibatches. Ce qui nous permet d'obtenir l'affichage suivant :



On remarque que l'évolution du temps d'exécution en fonction de la taille des minibatches semble relativement linéaire. Alors que la loss semble décroître de manière très rapide avec des tailles de minibatches assez petites (jusqu'à environ 50) puis ne plus beaucoup évoluer. On aurait donc tendance à choisir une taille de minibatch autour de 50.