
Évaluation de la performance de la classification d'images à l'aide de trois algorithmes

Lea-Marie Normandin, Zhibin Lu et Xiaocheng Liu

Introduction

Nous sommes intéressés par la classification d'images. Nous allons tester trois algorithmes pour en comparer leur performance et ainsi, choisir celui qui donne les meilleurs résultats. Pour chaque algorithme, le taux d'erreur sera calculé. Nous évaluerons pour différents noyaux et différentes valeurs d'hyper-paramètres afin de trouver le modèle optimal. Pour chaque algorithme sur chacune des bases de données, nous afficherons des exemples qui ont été difficilement classés et d'autres qui ont été facilement classés.

Les trois algorithmes utilisés seront le classificateur de Bayes naïf (un algorithme qui permet de classer les données en appliquant la règle de Bayes pour obtenir la probabilité à posteriori des classes aux points d'entrée (les composantes seront supposées indépendantes)), la machine à vecteurs de support (un algorithme qui apprend une fonction discriminante linéaire en maximisant la distance d'un hyper-plan à un exemple d'entraînement le plus proche de celui-ci), le réseau convolutif (un algorithme qui consiste en la construction d'un réseau de neurones convolutifs).

Nous utiliserons une base de données constituée de chiffres manuscrits (MNIST) ainsi qu'une base de données qui contient une grande quantité d'images en couleur de format 32x32 (CIFAR-10). Dans cette dernière, les images sont divisées en 10 classes et chaque classe contient 6000 images.

1 Classificateur de Bayes naïf

L'algorithme du classificateur de Bayes est construit autour de la règle de Bayes qui dit, dans notre contexte, que

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

où $X = (X_1, X_2, \dots, X_d) \in \mathbb{R}^d$ et $Y = c \in \{1, \dots, m\}$.

Le classificateur de Bayes naïf ajoute la condition d'indépendance entre les composantes du vecteur X étant donnée une certaine classe. On obtient donc que

$$P(X|Y = c) = P(X_1|Y = c)P(X_2|Y = c) \dots P(X_d|Y = c)$$

Dans les deux bases de données, il y a 10 classes.

1.1 MNIST

Un hyper-paramètre très souvent considéré dans l'algorithme du classificateur de Bayes est celui du lissage. Nous avons donc calculé les erreurs obtenues sur les ensembles d'entraînement et de validation pour différentes valeurs de cet hyper-paramètre. Ces erreurs sont affichées dans la table 1.

Table 1: Erreurs sur l'ensemble d'entraînement et l'ensemble de validation pour différentes valeurs de l'hyper-paramètre de lissage

Valeur de HP	Erreurs	
	Entraînement	Validation
0.2	0.15500	0.1673
0.5	0.15456	0.1642
2	0.15492	0.1646
20	0.15656	0.1658
40	0.15794	0.1668

30 Dans les deux cas, l'erreur la plus petite correspond à une valeur de l'hyper-paramètre de 0.5.

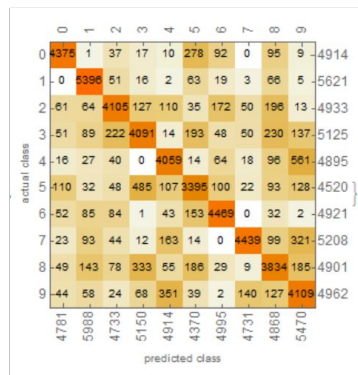


Figure 1: Matrice de confusion sur l'ensemble d'entraînement

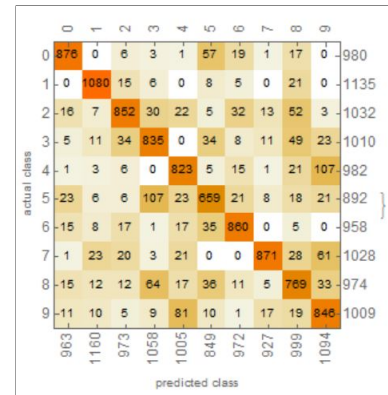


Figure 2: Matrice de confusion sur l'ensemble de test

31 Les matrices de confusion pour l'ensemble d'entraînement et l'ensemble de test avec un hyper-
32 paramètre de lissage de 0.5 sont affichées aux figures 1 et 2.

33 D'après ces dernières, les images où l'algorithme semble le plus confus sont celles représentant le
34 4. L'algorithme classe 561 images de 4 comme étant des 9 sur l'ensemble d'entraînement. Il y a
35 également le 5 qui est prédit comme étant dans la classe du 3 pour 485 images d'entrées sur ce même
36 ensemble. Pour l'ensemble de test, le 5 est souvent prédit comme un 3 et le 4 est souvent prédit par le
37 9 (107 fois chacun). Il demeure toutefois que la classe prédite la plus souvent correspond à la vraie
38 classe.

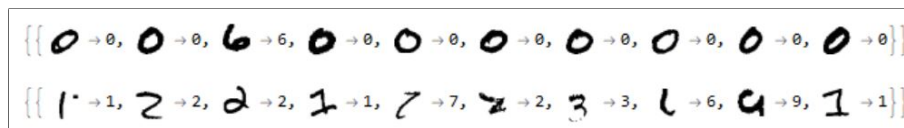


Figure 3: Exemples facilement classifiés et difficilement classifiés

39 On retrouve à la première ligne de la figure 3 quelques exemples de l'ensemble de test qui sont
40 facilement classifiés par l'algorithme et à la seconde ceux qui causent plus de difficultés.

41 1.2 CIFAR-10

42 Nous suivons les mêmes étapes sur cette nouvelle base de données.

43 L'hyper-paramètre de lissage conservé est 0.2 puisqu'il correspond à la plus petite erreur sur
44 l'ensemble d'entraînement (voir la table 2).

45 Les matrices de confusion sur l'ensemble d'entraînement et l'ensemble de test sont les figures 4
46 et 5. Pour l'ensemble d'entraînement, plusieurs cases à l'extérieur de la diagonale principale sont

Table 2: Erreurs sur l'ensemble d'entraînement et l'ensemble de validation pour différentes valeurs de l'hyper-paramètre de lissage

Valeur de HP	Erreurs	
	Entraînement	Validation
0.2	0.5879	0.7385
0.5	0.5881	0.7385
2	0.5882	0.7385
20	0.5884	0.7385
40	0.5893	0.7390

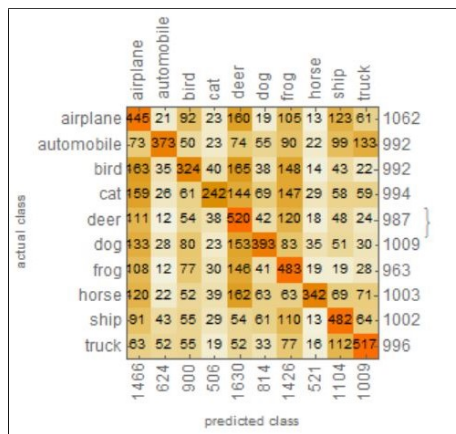


Figure 4: Matrice de confusion sur l'ensemble d'entraînement

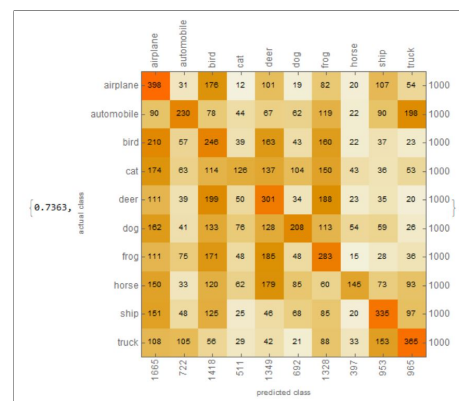


Figure 5: Matrice de confusion sur l'ensemble de test

assez foncées ce qui reflète grandement le taux d'erreur obtenu de 0.5879. Ce taux d'erreur, malgré qu'il soit le plus petit parmi les différentes valeurs pour l'hyper-paramètre, demeure un taux très grand. L'algorithme se trompe plus d'une fois sur deux. Pour l'ensemble de test, il est plus difficile de distinguer la diagonale principale qui est normalement un bon indicateur que l'algorithme a bien prédit chacune des classes. On voit ici que ce classificateur est très souvent mélangé et interprète mal une bonne quantité d'images. Cela peut fort probablement être expliqué par la supposition d'indépendance des composantes d'entrées et donc le manque de force de ce classificateur.

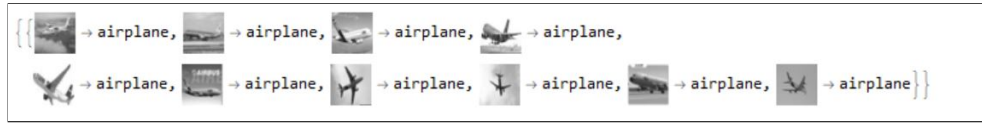


Figure 6: Exemples facilement classifiés

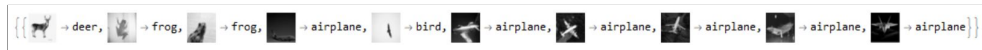


Figure 7: Exemples difficilement classifiés

Quelques exemples bien classifiés par l'algorithme se trouvent à la figure 6 et quelques exemples plus difficiles à classer se trouvent à la figure 7.

2 Machines à vecteurs de support

Les machines à vecteurs de support constituent un classificateur binaire. Les évaluations se font donc un à un ou un contre tous les autres. Un exemple de cette dernière option : soit l'exemple d'entrée fait partie d'une certaine classe c ou non. À partir de ce principe, l'algorithme apprend une fonction discriminante linéaire à marge maximale. Un hyperplan est également trouvé. Il est meilleur si il sépare parfaitement les données et si il se trouve à une distance éloignée du point d'entraînement le plus proche.

2.1 MNIST

Table 3: Erreurs sur l'ensemble d'entraînement et l'ensemble de test pour plusieurs noyaux

Type de noyau	Erreurs	
	Entraînement	Test
Linéaire	0.03832	0.0524
Polynomial	0.02696	0.0372
Gaussien	0.04654	0.0493

Dans le but de trouver un classificateur non-linéaire, nous commençons par comparer différents noyaux possibles. Nous trouvons le taux d'erreur sur l'ensemble d'entraînement et l'ensemble de test pour le noyau linéaire, le noyau polynomial et le noyau gaussien. Ces taux sont affichés à la table 3. Dans le logiciel utilisé, le noyau linéaire est d'équation $K(x_1, x_2) = x_1 x_2$, le noyau polynomial est d'équation $K(x_1, x_2) = \gamma(c + x_1 x_2)^d$ et le noyau gaussien est $K(x_1, x_2) = \exp\{\gamma|x_1 - x_2|^2\}$.

Le type de noyau qui obtient la plus petite erreur autant sur l'ensemble d'entraînement que sur l'ensemble de test est le noyau polynomial. Nous considérons donc différents degrés pour le polynôme (un hyper-paramètre) et en déterminons les erreurs sur l'ensemble de validation. Celles-ci sont affichées dans la table 4. L'erreur la plus petite correspond à un noyau de type polynomial de degré 9. Il reste donc finalement à choisir l'hyper-paramètre gamma γ .

Sur l'ensemble de validation, les erreurs pour différentes valeurs de gamma sont inscrites à la table 5. Étant donné que l'erreur la plus petite est à la fois pour un gamma de 1 et un gamma de 10, on conservera un gamma de 1.

Table 4: Erreurs sur l'ensemble de validation pour différents degrés du polynôme

Degré	Erreurs de validation
3	0.0397
4	0.0314
5	0.0257
9	0.0234

Table 5: Erreurs sur l'ensemble de validation pour différents gammas

γ	Erreurs de validation
0.0001	0.0578
0.001	0.0223
0.01	0.0183
0.1	0.0184
1	0.0180
10	0.0180

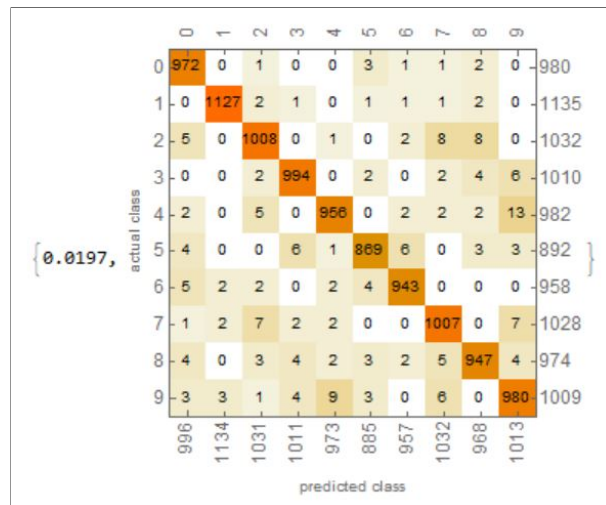


Figure 8: Matrice de confusion pour le modèle choisi

77 Sur l'ensemble de test, on obtient donc la matrice de confusion de la figure 8. Nous remarquons que
 78 les machines de vecteurs à support semblent mieux prédire les classes que le classificateur de Bayes
 79 naïf. En effet, en comparant les deux matrices de confusion (2 et 8) sur l'ensemble de test, il y a
 80 beaucoup moins de couleur foncée à l'extérieur de la diagonale principale pour les SVM que pour le
 81 classificateur de Bayes naïf.

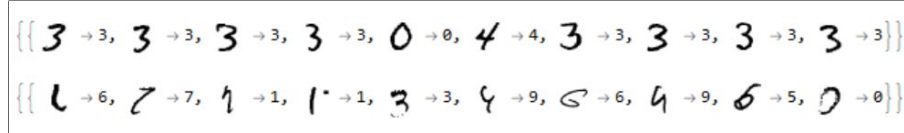


Figure 9: Exemples facilement classifiés et difficilement classifiés

82 Quelques exemples d'images facilement classifiées et difficilement classifiées se trouvent à la figure
 83 9.

Table 6: Erreurs sur l'ensemble d'entraînement et l'ensemble de test pour plusieurs noyaux

Type de noyau	Erreurs	
	Entraînement	Test
Linéaire	0.5469	0.7820
Polynomial	0.7524	0.7502
Gaussien	0.8272	0.8213

84 2.2 CIFAR-10

85 En suivant la même logique que pour la base de données MNIST, nous trouvons les erreurs pour les
 86 différents noyaux.

87 L'erreur de test la plus petite est celle du noyau de type polynomial comme le montre la table 6.

Table 7: Erreurs sur l'ensemble de validation pour différents degrés du polynôme

Degré	Erreurs de validation
3	0.75
4	0.845
5	0.7495
9	0.7375

88 En calculant les erreurs sur l'ensemble de validation, le degré 9 pour le polynôme est celui qui
 89 comporte l'erreur la plus petite. Ensuite, en ayant conservé un degré de 9, la valeur de l'hyper-
 90 paramètre gamma qui donne la plus petite erreur est 0.001 comme il est possible de voir à la table
 91 8.

92 La matrice de confusion sur l'ensemble de test est donc celle affichée à la figure 10 et les exemples
 93 facilement classifiés et difficilement classifiés sont aux figures 11 et 12 respectivement.

Table 8: Erreurs sur l'ensemble de validation pour différents gammas

γ	Erreurs de validation
0.0001	0.763
0.001	0.678
0.01	0.7155
1	0.776
10	0.7965

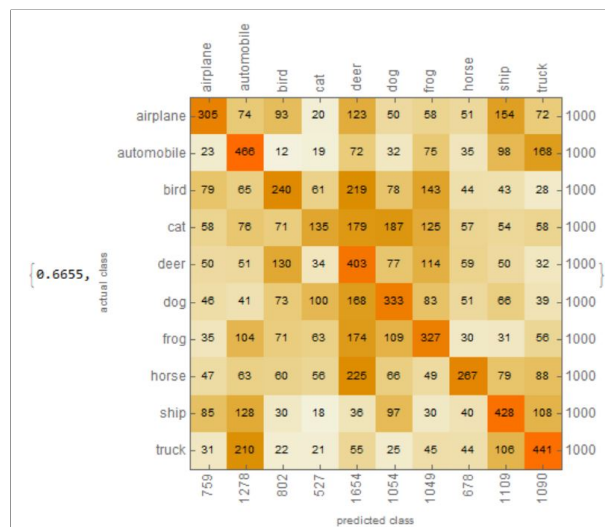


Figure 10: Matrice de confusion pour le modèle choisi



Figure 11: Exemples facilement classifiés

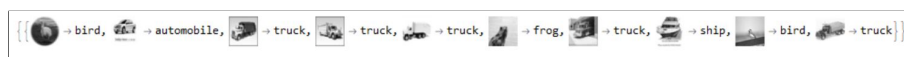


Figure 12: Exemples difficilement classifiés

94 3 Réseau convolutif

95 Le réseau de neurones convolutif considère une matrice de poids partagé (filtre) et fait une convolution
 96 sur l'image d'entrée dans le but d'extraire des informations spécifiques sans perdre de généralité. Il
 97 permet également de réduire le nombre de paramètres, c'est-à-dire le nombre de pixels.

98 3.1 MNIST

99 On commence par trouver les paramètres qui minimise le risque empirique. Cette étape, appelée
 100 entraînement, consiste en l'utilisation de la méthode de la descente du gradient pour l'optimisation
 101 des paramètres. À partir de cela, il est possible de varier certains hyper-paramètres en gardant les
 102 autres constants.

103 En considérant un nombre d'époques égal à 5, une taille des lots de 500 et un terme de régularisation
 104 nul, on fait varier le taux d'apprentissage.



Figure 13: Taux d'apprentissage = 5

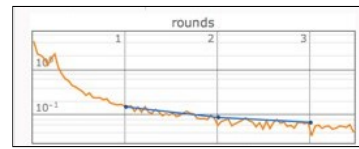


Figure 14: Taux d'apprentissage = 4

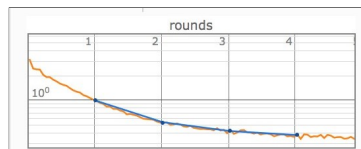


Figure 15: Taux d'apprentissage = 0.05

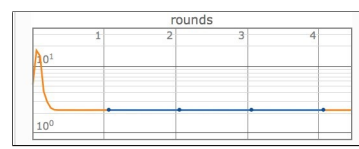


Figure 16: Taille des lots = 2000

105 Pour un taux d'apprentissage de 5, la méthode de descente du gradient n'aboutit pas. En effet, on
 106 remarque sur la figure 13 que le graphique affiche une courbe de loss constante à partir d'une certaine
 107 valeur.

108 Pour un taux d'apprentissage de 4 (voir la figure 14), malgré que la méthode de descente du gradient
 109 fonctionne, il semble y avoir de l'instabilité. On obtient, pour trois entraînements différents, des taux
 110 d'erreur très variables. Les trois taux d'erreur obtenus sont 0.0142, 0.0438 et 0.0133.

111 Pour un taux d'apprentissage de 0.05, la descente du gradient est lente. Cela est remarqué en regardant
 112 l'échelle du graphique de la figure 15.

113 Maintenant, en considérant un nombre d'époques égal à 5, un taux d'apprentissage de 3 et un terme
 114 de régularisation nul, on fait varier la taille des lots.

115 Pour une taille des lots de 2000, on remarque qu'une trop grosse taille des lots empêche l'utilisation
 116 de la descente du gradient. Après plusieurs époques, on voit à la figure 16 que la courbe est constante.

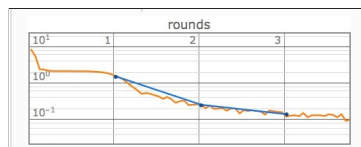


Figure 17: Taille des lots = 1000

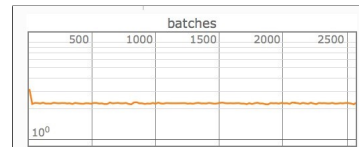


Figure 18: Taille des lots = 5

117 Pour une taille des lots de 1000, divers entraînements arrivent à différents taux d'erreur. On obtient
 118 les taux suivants : 0.0297, 0.0657 et 0.0205. On a donc, une fois de plus, instabilité comme le montre
 119 la figure 17.

120 Pour une taille des lots de 5, la descente du gradient est très lente pour ne pas dire inexistante (voir la
 121 figure 18).

122 Ensuite, on ajuste l'hyper-paramètre afin de contrôler le terme de régularisation.
 123 En considérant un nombre d'époques de 30, un taux d'apprentissage de 1.8 et une taille des lots de
 124 100, on tente de trouver la valeur de cet hyper-paramètres qui donne la meilleure performance.

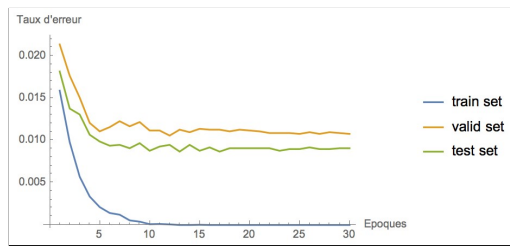


Figure 19: Terme de régularisation = 0.000000000001

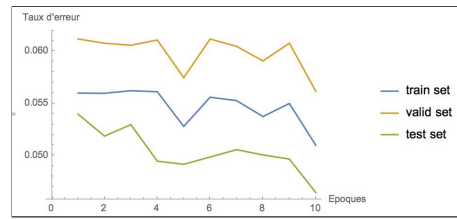


Figure 20: Terme de régularisation = 0.05

125 Pour un terme de régularisation nettement inférieur à zéro (0.000000000001), on remarque dans
 126 la figure 19 une situation de sur-apprentissage. En effet, on remarque une différence significative
 127 entre l'erreur sur l'ensemble d'entraînement et celle sur l'ensemble de test plus le nombre d'époques
 128 augmentent. Après l'époque 10, on constate que le taux d'erreur sur l'ensemble d'entraînement
 129 est proche de 0 alors que celui pour les deux autres ensembles est entre 0.010 et 0.014. Ce taux
 130 d'erreur sur l'ensemble de test est tout de même petit ce qui est expliqué, entre autres, par la taille de
 131 l'échantillon.

132 Pour un terme de régularisation de 0.05, on remarque une situation de sous-apprentissage. La figure
 133 20 illustre effectivement un biais assez grand, mais une petite variance.

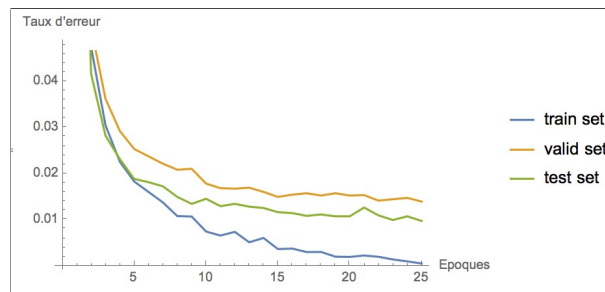
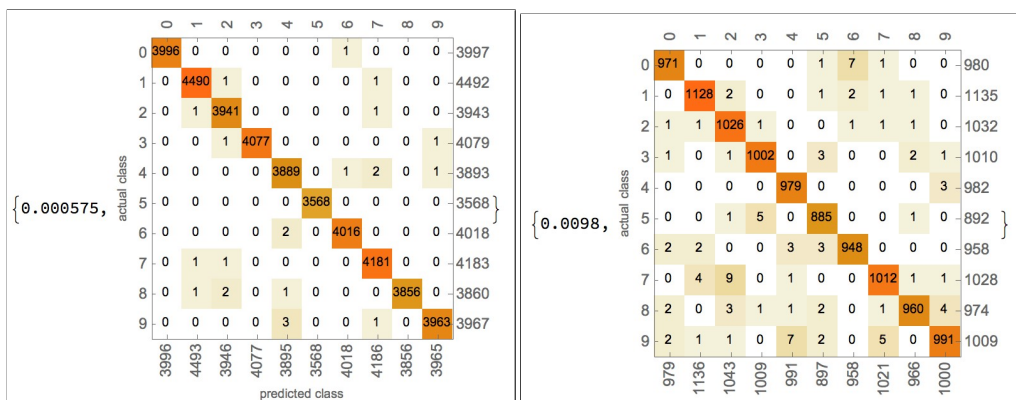


Figure 21: Graphique obtenu avec les paramètres et hyper-paramètres choisis

134 Suite à cette évaluation, nous avons décidé de fixer le nombre d'époques à 25, le taux d'apprentissage
 135 à 1.7, la taille des lots à 240 et le terme de régularisation à 0.000001. Pour ces valeurs, on obtient le
 136 graphique du taux d'erreur en fonction des époques à la figure 21.

137 On obtient donc les taux d'erreur et les matrices de confusion sur les ensembles d'entraînement et de
 138 test, respectivement aux figures 22 et 23.

139 Pour évaluer le niveau de la performance de l'algorithme, nous considérons l'entropie, c'est-à-dire
 140 le niveau d'incertitude par rapport aux différentes images d'entrée. Une entropie élevée signifie la
 141 possibilité de confusion et donc aussi la possibilité d'échec de l'algorithme. Une entropie faible
 142 signifie, pour sa part, une bonne chance de réussite de l'algorithme. À la figure 24, la première
 143 ligne contient les exemples où l'algorithme échoue alors que la seconde contient les exemples
 144 où l'algorithme réussit le mieux. D'après les matrices de confusion, nous pouvons constater que
 145 l'algorithme semble parfois confondre le 6 avec le 0, le 7 avec le 2 et le 4 avec le 9.



{ 2, 4, 9, 7, 0, 7, 8, 6, 5, 1, 2, 9, 8, 0, 1, 5, 1, 0, 0, 8, 9, 9, 0, 2, 9, 1, 7, 8, 2, ; }

{ 3, 3, 3, 4, 3, 6, 6, 6, 3, 3, 3, 3, 3, 4, 3, 3, 6, 4, 3, 3, 4, 6, 3, 2, 3, 3, 3, 6, 6, 3 }

Figure 24: Exemples

146 3.2 CIFAR-10

147 Suite à l'expérience sur la base de données MNIST, on considère un nombre d'époques entre 10 et 20,
 148 une taille des lots entre 200 et 300, un terme de régularisation de 0.0000001, un taux d'apprentissage
 149 de 2.4 et deux couches cachées contenant respectivement 800 et 100 neurones.

150 On commence par vérifier l'influence de certains hyper-paramètres tels que le nombre de couches de
 151 convolution, le nombre de couches de pooling ainsi que leurs dimensions respectives.

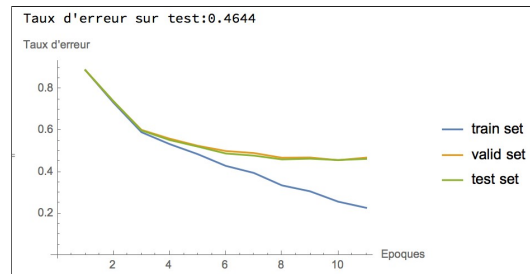


Figure 25: Premier test

152 On commence à tester l'algorithme dans une situation où l'on a une couche de convolution et une
 153 couche de pooling. La couche de convolution contient 30 filtres de dimension 5x5 et la couche de
 154 pooling est composée d'un seul filtre de dimension 2x2. La figure 25 représente le graphique du taux
 155 d'erreur en fonction des époques.

156 On constate une situation de sous-apprentissage. En effet, la capacité du modèle n'est pas suffisamment
 157 élevée.

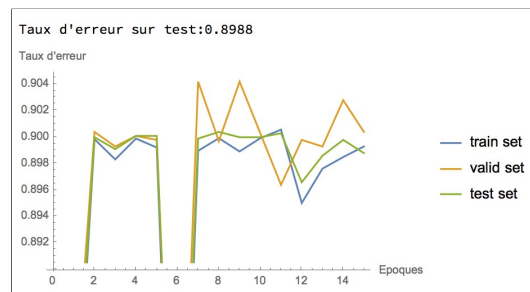


Figure 26: Deuxième test

158 On tente maintenant de considérer deux couches de convolution et deux couches de pooling. La
 159 première couche de convolution est composée d'une trentaine de filtres de dimension 9x9 et la
 160 seconde est composée de 90 filtres de même taille. Les deux couches de pooling maintiennent une
 161 dimension de 2x2 avec un seul filtre chacun. La figure 26 affiche un résultat incohérent qui nous
 162 permet de dire que la taille des filtres est trop grande pour des images de format 32x32.

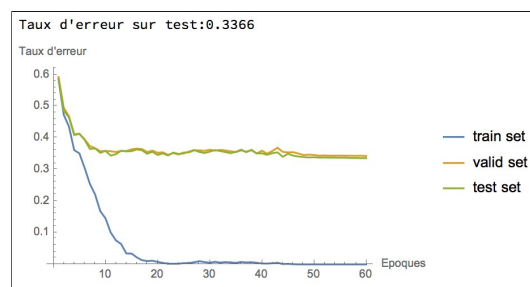


Figure 27: Troisième test

163 On recommence donc la dernière expérience en ne changeant que la taille des filtres des couches de
 164 convolution. La taille est établie à 5x5. La figure 27 affiche un résultat meilleur que ceux des deux
 165 tests précédents. Le taux d'erreur est plus bas.

166 Cette amélioration marquée nous pousse à augmenter davantage le nombre de couches et à diminuer la
 167 taille des filtres. Pour trois couches de convolution avec des filtres de dimension 3x3 et trois couches
 168 de pooling avec un filtre de 2x2 chacun. Le nombre de filtres pour chaque couche de convolution est
 169 60, 120 et 240 respectivement.

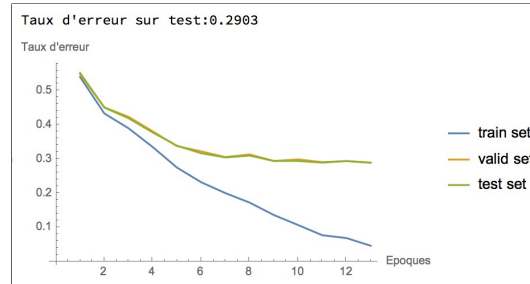


Figure 28: Quatrième test

170 On obtient donc le graphique de la figure 28.

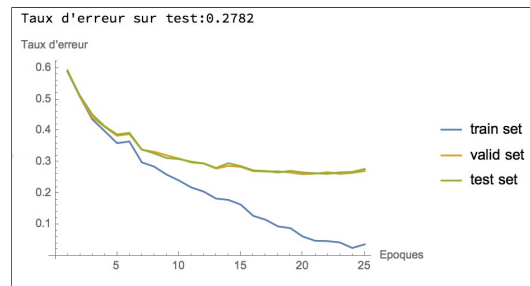


Figure 29: Cinquième test

171 Jusqu'à maintenant, la fonction exécutée sur chacune des couches de pooling est la fonction max. En
 172 considérant plutôt la fonction max seulement sur la première couche et la fonction moyenne sur les
 173 deux autres, le résultat obtenu est encore meilleur que les précédents comme il est possible de voir à
 174 la figure 29. Le taux d'erreur est à son plus bas.

actual class	predicted class										
	ship	airplane	automobile	bird	cat	deer	dog	frog	horse	truck	
ship	827	39	38	25	26	8	16	10	4	7	1000
airplane	77	895	30	74	29	37	12	18	13	15	1000
automobile	28	11	868	13	13	3	13	22	3	28	1000
bird	7	37	6	642	55	87	64	75	23	4	1000
cat	14	11	7	80	548	55	186	71	20	8	1000
deer	11	10	3	79	58	879	53	55	49	3	1000
dog	4	7	0	49	139	44	885	35	26	1	1000
frog	2	1	4	47	52	25	31	833	4	1	1000
horse	3	6	8	32	55	68	94	8	720	6	1000
truck	30	26	108	17	38	8	20	21	19	713	1000
	1003	843	1070	1058	1013	1014	1184	881	786		

Figure 30: Matrice de confusion CIFAR-10

175 La matrice de confusion obtenue en exécutant le dernier algorithme choisi est celle affichée à la figure
 176 30.



Figure 31: Entropie plus élevée



Figure 32: Entropie plus basse

177 Les images ayant les entropies les plus hautes et les plus basses sont aux figures 31 et 32 respective-
178 ment.

179 On constate, à partir de la matrice de confusion que 108 camions ont été considérés comme des
180 automobiles, que 186 chats ont été classés comme chien par l'algorithme et que 139 chiens ont été
181 placés dans la catégorie chat. Les images dont l'entropie est la plus élevée semblent avoir moins de
182 contrastes que celles dont l'entropie est plus basse. Cela explique en grande partie la confusion de
183 l'algorithme. Les images les mieux classées sont les automobiles et les chevaux.

184 4 Conclusion

185 Nous avons remarqué que le taux d'erreur décroît dans l'ordre des algorithmes présentés pour la
186 base de données MNIST. L'algorithme du réseau convolutif serait donc le meilleur. Pour la base de
187 données CIFAR-10, les deux premiers classificateurs ne donnent pas un super résultat et le réseau
188 convolutif est de loin meilleur aux précédents. Le réseau convolutif est donc un bon classificateur
189 pour un problème avec des images. Pour un projet futur, il serait intéressant d'essayer d'autres
190 algorithmes sur ces mêmes bases de données afin de voir s'il y en a d'autres qui donnent de meilleurs
191 résultats.

192 Répartition

193 La répartition fût équitable entre les membres de l'équipe.