

Rapport de séance semaine du 17-12-2018

Cette semaine je me suis concentré sur une seule chose, réussir à importer les données de la base de l'EMNIST dans deux tableau python.

Cette étape semble facile, mais ce n'est pas du tout le cas. Il y a très peu de documentation sur l'EMNIST car la plupart des exemples sont fait avec l'MNIST (une version bien moins complète).

La base de l'EMNIST est une ressource divisée en plusieurs fichiers et qui existe en plusieurs versions, d'un poids de 700 mo, elle ne comporte pas moins de 700 000 caractères (chiffres et lettres).

On le charge sous forme de deux tableau à 3 dimensions :

Le premier contient l'index du caractère et les pixels.

Le deuxième contient le caractère écrit en version ordinateur.

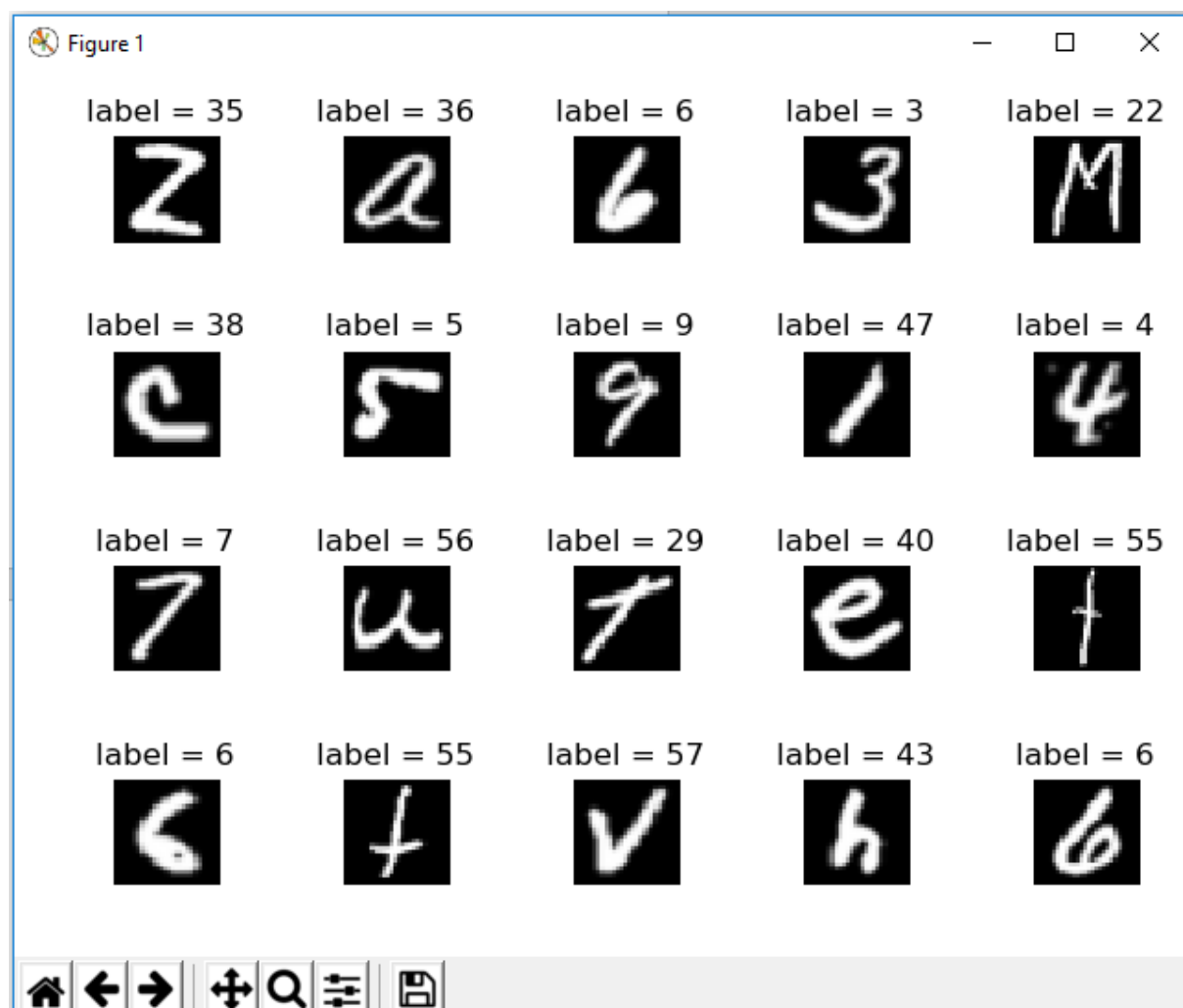
Au bout d'environ 3 heures, j'ai réussi à mixer plusieurs codes pour pouvoir charger les données (il faut environ 1 minutes à un ordinateur bien au-dessus de la moyenne pour pouvoir les charger (exécution en single core), et nécessite au minimum 5.5 go de mémoire vive non utilisé).

Le code est au final minuscule car il j'ai trouvé une classe (MNIST) déjà toute faite qui s'occupe de tout, sinon le code aurait était très compliqué (manipulation de gros tableaux avec beaucoup de transformations).

J'ai écrit ce tableau alliant label et correspondance

0	0	30	U
1	1	31	V
2	2	32	W
3	3	33	X
4	4	34	Y
5	5	35	Z
6	6	36	a
7	7	37	b
8	8	38	c
9	9	39	d
10	A	40	e
11	B	41	f
12	C	42	g
13	D	43	h
14	E	45	i
15	F	46	j
16	G	47	k
17	H	48	l
18	I	49	m
19	J	50	o
20	K	51	p
21	L	52	q
22	M	53	r
23	N	54	s
24	O	55	t
25	P	56	u
26	Q	57	v
27	R	58	w
28	S	59	x
29	T	60	y
		61	z

Après avoir finalement créer le tableau, j'ai tester, et voici quelques résultats (une infime partie) :



On peut maintenant utiliser ces données pour créer le premier réseau de neurones capable de reconnaître les caractères .

Information : Le fichier de données est trop lourd pour être portée sur GitHub