

# Conversion de données vers MNIST

Vivien Kraus

29 mars 2017

## Table des matières

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Compilation</b>	<b>1</b>
<b>3</b>	<b>Implémentation</b>	<b>2</b>
3.1	Composition du fichier de sortie . . . . .	2
3.2	Indexation . . . . .	3
3.3	Lecture des images . . . . .	5
3.4	Fabrication du fichier à partir de la base de données . . . . .	6
3.5	Application principale . . . . .	9
3.6	Compilateur . . . . .	9
<b>4</b>	<b>Tangle</b>	<b>9</b>

## 1 Introduction

Le but de ce programme est de convertir des bases de données de *deep learning* vers le format MNIST, tel que défini sur le site de Yann LeCun<sup>1</sup>.

Ici, nous allons avoir pour les données d'entrée des blocs de matrices, donc le « nombre magique » IDX est 4.

Nous allons charger des données sous la forme suivante : la base de données d'images est séparée entre **tests** et **apprentissage**, et chaque dossier contient des images dont le nom est **target\_identifiant\_vue.jpg**, avec **target** le nom de la conclusion en toutes lettres, **identifiant** un nombre qui définit un individu, et **vue** le numéro de la vue. Les identifiants sont discontinus.

Les noms des cibles sont transformés en numéro, pour retrouver les noms, nous créons aussi un fichier d'indice avec une conclusion par ligne.

## 2 Compilation

Vous pouvez télécharger l'archive source ici. Décompressez-la. Vous devez avoir sur votre ordinateur :

---

1. <http://yann.lecun.com/exdb/mnist/>

- `ocamlfind`;
- `ocamlbuild`;
- `camlimages` (avec le module `jpeg`).

Le meilleur moyen est d'utiliser OPAM, plus de détails ici.

Une fois l'archive téléchargée et dépaquetée, il faut compiler. Si vous avez enregistré dans la variable `PATH_TO_SHREC_DATABASE` le chemin d'accès vers la base SHREC (un dossier contenant deux sous-dossiers, **test** et **train**), et dans `MNIST_OUTPUT_DIR` le chemin où enregistrer les fichiers `idx`, vous pouvez utiliser la commande suivante :

```
cd idx
ocamlbuild -use-ocamlfind \
    src/main.native -- \
    $PATH_TO_SHREC_DATABASE $MNIST_OUTPUT_DIR
```

Listing 1: Commande pour exécuter le programme

## 3 Implémentation

Nous allons utiliser Objective CAML.

### 3.1 Composition du fichier de sortie

Le but est de calculer la composition du fichier de sortie. On obtient l'interface suivante :

```
1  (** Compute the file content from the input data. *)
2
3  (** [write_idx dimensions write f] will call [write] in turn to all
4  bytes. The input is given by [f coordinates], where [coordinates] are
5  the coordinates, such that the i-th coordinate is between 0 and the
6  i-th dimension. *)
7  val write_idx: int list -> (int -> unit) -> ((int list -> int) -> unit)
```

Et voici l'implémentation :

```
1  let rec to_net acc = function
2    | n when n <= 0 -> acc
3    | n -> to_net ((n mod 256) :: acc) (n / 256)
4
5  let int32_to_net x =
6    match to_net [] (Int32.to_int x) with
7    | [] -> [0; 0; 0; 0]
8    | [x] -> [0; 0; 0; x]
9    | [x; y] -> [0; 0; x; y]
10   | [x; y; z] -> [0; x; y; z]
```

```

11 | [x; y; z; t] -> [|x; y; z; t|]
12 | _ -> failwith "Impossible"
13
14 let write_idx dimensions write_char =
15   let write_uint n =
16     let n = Int32.of_int n in
17     let n = int32_to_net n in
18     Array.iter write_char n
19   in
20   let rec range_aux read_char coordinates = function
21     | [] -> write_char (read_char coordinates)
22     | n :: tl ->
23       for i=0 to n-1 do
24         range_aux read_char (i :: coordinates) tl
25       done
26   in
27   let range dimensions read_char = range_aux read_char [] dimensions in
28   write_char 0;
29   write_char 0;
30   write_char 8;
31   write_char (List.length dimensions);
32   List.iter write_uint dimensions;
33   range dimensions

```

### 3.2 Indexation

Ici, nous parcourons la base de données SHREC pour enregistrer les cibles. Nous obtenons une interface un peu plus complexe :

```

1  (** Search both train and test databases *)
2
3  (** Keep the association between target names and values *)
4  type database =
5    {
6      (** Database directory name *)
7      dirname: Bytes.t;
8      (** list of targets *)
9      targets: Bytes.t array;
10     (** Number of views *)
11     n_views: int;
12     (** List of train items - unsorted, conclusion * item ID *)
13     train: (int * Bytes.t) array;
14     (** List of test items *)
15     test: (int * Bytes.t) array;
16   }

```

```

17
18 (** Read the database under given directory *)
19 val read_database: Bytes.t -> database

```

Et une implémentation qui utilise le module Sys (<https://caml.inria.fr/pub/docs/manual-ocaml/libref/Sys.html>) :

```

1  type database =
2      {
3          dirname: Bytes.t;
4          targets: Bytes.t array;
5          n_views: int;
6          train: (int * Bytes.t) array;
7          test: (int * Bytes.t) array;
8      }
9
10 let randomize t =
11     let swap i j =
12         let x = t.(i) in
13         let () = t.(i) <- t.(j) in
14         t.(j) <- x
15     in
16     for i=1 to Array.length t - 1 do
17         swap i (Random.int i)
18     done
19
20 let read_database path =
21     let n_views = ref 0 in
22     let n_targets = ref 0 in
23     let targets = Hashtbl.create 100 in
24     let rec get_target name =
25         try
26             Hashtbl.find targets name
27         with Not_found ->
28             let () = Hashtbl.add targets name !n_targets in
29             let () = incr n_targets in
30             get_target name
31     in
32     let analyze_path path =
33         let b = Filename.basename path in
34         let triple = Filename.chop_extension b in
35         try
36             let sep_1 = Bytes.index triple '_' in
37             let sep_2 = Bytes.index_from triple (sep_1 + 1) '_' in
38             let start_target = 0 in

```

```

39   let end_target = sep_1 in
40   let len_target = end_target - start_target in
41   let start_id = sep_1 + 1 in
42   let end_id = sep_2 in
43   let len_id = end_id - start_id in
44   let start_view = sep_2 + 1 in
45   let end_view = Bytes.length triple in
46   let len_view = end_view - start_view in
47   let (t, i, v) =
48     (Bytes.sub triple start_target len_target,
49      Bytes.sub triple start_id len_id,
50      Bytes.sub triple start_view len_view) in
51   let v = int_of_string v in
52   let t = get_target t in
53   let () = if v > !n_views then n_views := v in
54   if v = 1 then
55     [(t, i)]
56   else []
57   with _ -> []
58 in
59 let analyze_dir path =
60   let files = Sys.readdir path in
61   let () = randomize files in
62   let files = Array.to_list files in
63   let files = List.map (fun f -> Filename.concat path f) files in
64   let items = List.map analyze_path files in
65   List.concat items
66 in
67 let train_items = analyze_dir (Filename.concat path "train") in
68 let test_items = analyze_dir (Filename.concat path "test") in
69 let all_targets = Array.make !n_targets "" in
70 let () = Hashtbl.iter (fun name i -> all_targets.(i) <- name) targets in
71 {
72   dirname = path;
73   targets = all_targets;
74   n_views = !n_views;
75   train = Array.of_list train_items;
76   test = Array.of_list test_items;
77 }

```

### 3.3 Lecture des images

Ici, on utilise la bibliothèque Camlimages ([gallium.inria.fr/camlimages/](http://gallium.inria.fr/camlimages/), documentation en ligne [file:///usr/share/doc/libcamlimages-ocaml-doc/html/api/index.html](http://file:///usr/share/doc/libcamlimages-ocaml-doc/html/api/index.html)).

L'interface est très simple, on récupère la matrice des pixels :

```

1  (** Read images *)
2
3  (** Read the image under the given filename *)
4  val read_image: Bytes.t -> int array array

```

L'implémentation est extrêmement simple :

```

1  let read_image filename =
2    let image = Images.load filename [] in
3    let (n, p, get_rgb) =
4      match image with
5      | Images.Index8 i8t ->
6        (i8t.Index8.height, i8t.Index8.width,
7         fun i j -> Index8.get_rgb i8t j i)
8      | Images.Rgb24 rgb24t ->
9        (rgb24t.Rgb24.height, rgb24t.Rgb24.width,
10         fun i j -> Rgb24.get rgb24t j i)
11     | Images.Index16 i16t ->
12       (i16t.Index16.height, i16t.Index16.width,
13        fun i j -> Index16.get_rgb i16t j i)
14     | Images.Rgba32 rgba32t ->
15       (rgba32t.Rgba32.height, rgba32t.Rgba32.width,
16        fun i j ->
17          let rgba = Rgba32.get rgba32t j i in
18          rgba.Color.color)
19     | Images.Cmyk32 cmyk32t ->
20       failwith "Error: this is a CMYK image, I can't handle it"
21   in
22   let get i j =
23     let rgb = get_rgb i j in
24     (rgb.Color.r + rgb.Color.g + rgb.Color.b) / 3
25   in
26   Array.init n (fun i -> Array.init p (get i))

```

### 3.4 Fabrication du fichier à partir de la base de données

On effectue ici la conversion à proprement parler. L'interface fournit une fonction qui prend en entrée l'indexation et le dossier où écrire les tables idx :

```

1  (** Convert image database to idx *)
2
3  (** Convert the database to four files in the given directory. *)
4  val convert: Index.database -> Bytes.t -> unit

```

L'implémentation est un peu plus longue ; il faut charger les images une par une :

```
1 let convert db dirname =
2   let open Index in
3   let open Compute in
4   let load_view dirname (target, item_id) view =
5     let basename =
6       Printf.sprintf
7         "%s_%s_%03d.jpg"
8         (db.targets.(target))
9         item_id
10        (view + 1) in
11   let name = Filename.concat
12     (Filename.concat db.dirname dirname)
13     basename in
14   Read_images.read_image name
15 in
16 let do_database write_dir items =
17   let i_item = ref 0 in
18   let i_view = ref 0 in
19   let i = ref 0 in
20   let j = ref 0 in
21   let view = ref (load_view dir items.(0) !i_view) in
22   let get () =
23     let v = !view in
24     let i = !i in
25     let j = !j in
26     v.(i).(j)
27   in
28   let next () =
29     let () = incr j in
30     let v = !view in
31     if !j >= Array.length (v.(!i)) then
32       let () = incr i in
33       let () = j := 0 in
34       if !i >= Array.length v then
35         let () = incr i_view in
36         let () = i := 0 in
37         let () = j := 0 in
38         let () =
39           if !i_view >= db.n_views then
40             let () = incr i_item in
41             let () = i_view := 0 in
42             ()
43         in
44         if !i_item < Array.length items then
```

```

45         view := load_view dir items.(!i_item) !i_view
46     in
47     let get_next coord =
48         let x = get () in
49         let () = next () in
50         x
51     in
52     let dimensions =
53         [Array.length items;
54          db.n_views;
55          Array.length !view;
56          Array.length !view.(0)] in
57     write_idx dimensions write get_next
58 in
59 let do_labels write items =
60     let i = ref 0 in
61     let get_next _ =
62         let (x, _) = items.(!i) in
63         let () = incr i in
64         x
65     in
66     write_idx [Array.length items] write get_next
67 in
68 let train_images = open_out (Filename.concat
69                             dirname "train-images-idx4.ubyte") in
70 let train_labels = open_out (Filename.concat
71                             dirname "train-labels-idx1.ubyte") in
72 let test_images = open_out (Filename.concat
73                             dirname "test-images-idx4.ubyte") in
74 let test_labels = open_out (Filename.concat
75                             dirname "test-labels-idx1.ubyte") in
76 let dict = open_out (Filename.concat dirname "targets.txt") in
77 let () =
78     for i=0 to Array.length db.targets - 1 do
79         output_string dict (db.targets.(i) ^ "\n")
80     done
81 in
82 let write_stream c =
83     let c = char_of_int c in
84     output_char stream c
85 in
86 do_database (write train_images) "train" db.train;
87 do_database (write test_images) "test" db.test;
88 do_labels (write train_labels) db.train;
89 do_labels (write test_labels) db.test

```



### 3.5 Application principale

On utilise les deux arguments du programme comme chemin d'accès aux dossiers de la base SHREC et MNIST :

```
1 match Sys.argv with
2 | [|_; input; output|] ->
3     let db = Index.read_database input in
4     Convert.convert db output
5 | _ ->
6     Printf.printf "Usage: shrec-to-mnist [SHREC database] [MNIST output dir]\n%!"
```

### 3.6 Compilateur

Pour compiler le programme, on utilise `ocamlbuild`. Voir le manuel en ligne.

Nous voulons simplement compiler et linker les paquets Ocamlfind `camlimages` et `camlimages.jpeg`. Nous avons donc besoin d'un fichier `_tags` pour instruire ces dépendances externes.

```
1 <src/read_images.ml>: package(camlimages), package(camlimages.jpeg)
2 <src/main.{byte,native}>: package(camlimages), package(camlimages.jpeg)
```

Nous utilisons les tags `package()`, ce qui signifie qu'il faut lier le plugin `ocamlfind` au plugin de compilation ; pour ce faire on passe `-with-ocamlfind` à la compilation.

## 4 Tangle

Ici, nous extrayons le code source.

```
(org-babel-tangle)
```

Listing 2: Extraire le code source

```
tar cf shrec-to-mnist.tar.gz \  
  idx/src/compute.ml \  
  idx/src/compute.mli \  
  idx/src/index.ml \  
  idx/src/index.mli \  
  idx/src/read_images.ml \  
  idx/src/read_images.mli \  
  idx/src/convert.ml \  
  idx/src/convert.mli \  
  idx/src/main.ml \  
  idx/_tags
```

Listing 3: Empaquetage

L'archive est donc disponible [ici](#).