

Table de hachage : implémentation

INFORMATIQUE COMMUNE - TP n° 3.1 - Olivier Reynet

À la fin de ce chapitre, je sais :

- ☞ utiliser les listes Python
- ☞ écrire des fonctions en Python
- ☞ utiliser une bibliothèque en l'important correctement
- ☞ expliquer le fonctionnement d'une table de hachage

L'objectif de ce TP est de construire une table de hachage «à la main», un équivalent des `dict` Python. Dans ce but, il faut dans un premier temps disposer d'une fonction de hachage adaptée. C'est l'objet de la première partie.

On rappelle sur la figure 1 le principe du dictionnaire (ou table de hachage).

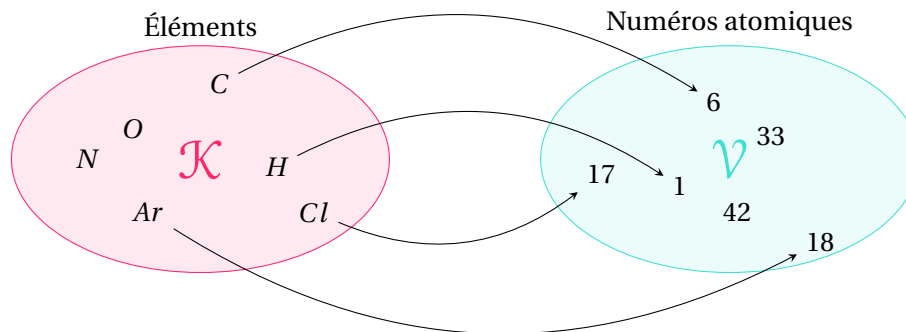


FIGURE 1 – Illustration du concept de dictionnaire : tableau associatif reliant une clef à un numéro atomique.

A Fonctions de hachage et uniformité

Pour être adaptée à l'usage par une table de hachage, une fonction de hachage doit être :

1. rapide,
2. cohérente : pour une même clef, on obtient un même code,
3. injective : pour des clefs différentes, on obtient des codes différents. Pour des codes identiques, les clefs sont nécessairement identiques. Dans le cas contraire, on obtient un **collision** qu'on cherche à minimiser.

- uniformément répartie : pour des clefs qui se ressemblent, les codes obtenus doivent être très différents, ceci pour limiter les collisions.

Une distribution uniforme des clefs dans l'espace d'arrivée peut être obtenue en utilisant des générateurs aléatoires. Les générateurs à congruence linéaires, c'est à dire les fonctions du type $(ax + b) \bmod n$ sont de bons candidats pour les fonctions de hachage, pourvu qu'on choisisse bien les constantes a et b du générateur.

Dans un premier temps, on opère l'**encodage de la chaîne de caractère en nombre entier** $\gamma(s)$ pour chaque chaîne de caractère s de la manière suivante :

$$\gamma(s) = \sum_{k=0}^{|s|-1} \text{ascii}(s_k) 2^{k+8} \quad (1)$$

où $\text{ascii}(s_k)$ est le code ASCII associé au caractère d'indice k de s . Cette étape est importante, car elle permet déjà de générer des codes différents pour des chaînes similaires.

P En Python, la fonction `ord` permet d'obtenir le code ASCII associé à un caractère (`documentation`).

Dans un second temps, on cherche à **compresser la valeur encodée dans l'intervalle des index possibles** $\llbracket 0, n \rrbracket$. Si n est la taille de la table de hachage, on peut choisir :

- d'utiliser simplement une division :

$$h_d : (s, n) \rightarrow \gamma(s) \bmod n \quad (2)$$

- d'utiliser une multiplication et une division :

$$h_\alpha : (s, n) \rightarrow \lfloor n \times (\alpha \gamma(s) \bmod 1) \rfloor \quad (3)$$

$\alpha \in]0, 1[$ étant une constante réelle.

Le choix d'une fonction de hachage est délicat et il n'existe pas de méthode pour atteindre l'optimal.

A1. Expliquer pourquoi la fonction γ renvoie un nombre unique associé à une chaîne de caractères.

Solution : Comme les caractères ASCII sont codés sur huit bits au maximum et que la fonction γ décale de $k+8$ bits vers la gauche chaque valeur $\text{ascii}(s_k)$, alors le nombre obtenu dépend des codes des lettres $\text{ascii}(s_k)$ et de leur position dans le mot (k). Les mots proches comportant les mêmes lettres mais pas dans le même ordre ne produisent donc pas le même code et les codes sont même distants les uns des autres :

```
1 print(gamma("abaa"), gamma("aaba"), gamma("aaab"), gamma("baaa"))
2 #372992 373504 374528 372736
3 print(gamma("choir"), gamma("music"), gamma("piano"), gamma("song"))
4 #874240 826112 861696 409856
```

A2. Coder les fonctions γ , h_d et h_α en Python. Les paramètres n et α pourront être pris par défaut à 47057 et $\frac{\sqrt{5}-1}{2}$. Ne pas oublier également que les puissances de deux peuvent être facilement calculées grâce aux opérateurs de décalage binaire.

Solution :

```

1
2 TABLE_SIZE = 47057
3 ALPHA = (1 + 5 ** 0.5) / 2
4
5 def gamma(s):
6     g = 0
7     for k in range(len(s)):
8         g += ord(s[k]) << (k + 8)
9     return g
10
11
12 def hd(key, table_size=TABLE_SIZE):
13     return gamma(key) % table_size
14
15
16 def hm(key, table_size=TABLE_SIZE, alpha=ALPHA):
17     return math.floor(table_size * (alpha * gamma(key) % 1))

```

- A3. Importer tous les mots contenus dans le fichier `"english_words.csv"` dans une liste.

On cherche à tester l'uniformité de la distribution des codes obtenues des fonctions de hachage. On peut facilement vérifier ceci en utilisant le test de Kolmogorov-Smirnov et la bibliothèque `scipy` et l'instruction :

```

1     scipy.stats.kstest(codes, "uniform")
2     #KstestResult(statistic=0.0012179563749926126, pvalue=0.49280753163611735)

```

Si le paramètre `p_value` est plus grand que 0.05, alors la distribution peut être considérée comme uniforme. Le paramètre `statistic` donne une mesure de la distance entre les deux distributions.

- A4. Écrire une fonction dont le prototype est `uniform_test(h, table_size)` dont le paramètre `h` est une fonction de hachage et `table_size` la taille de la table de hachage. Cette fonction renvoie le résultat du test de Kolmogorov-Smirnov entre une distribution uniforme et la distribution des codes obtenus avec `h` sur l'ensemble des mots du fichier `"english_words"`. La fonction de `scipy` nécessite un tableau d'entrée Numpy dont les données sont de type `float`.

Solution :

```

1 def uniformity_test(h, table_size):
2     codes = []
3     f = open("english_words.txt", "r")
4     for line in f:
5         words = line.split('\n')
6         hash_code = h(words[0], table_size)
7         codes.append(hash_code)
8     f.close()
9     codes = np.array(codes, dtype=float)
10    codes = codes / np.max(codes)
11    return scipy.stats.kstest(codes, "uniform")

```

- A5. Observer les résultats de la fonction précédente pour h_d et h_a en faisant varier la taille de la table de hachage. Que pouvez-vous en conclure?

Solution : La fonction h_d est sensible au choix de la taille de la table de hachage alors que h_a ne l'est pas.

- A6. Afficher les histogrammes associés aux différentes distributions de codes obtenues à l'aide de la bibliothèque `matplotlib` et à la fonction `hist`.

Solution :

```
1  def plot_hist(h):
2      codes = []
3      f = open("english_words.txt", "r")
4      for line in f:
5          words = line.split('\n')
6          hash_code = h(words[0], TABLE_SIZE)
7          codes.append(hash_code)
8      f.close()
9      codes = np.array(codes, dtype=float)
10     codes = codes / np.max(codes)
11     plt.hist(codes, 50) # range=(np.min(codes), np.max(codes))
12     plt.title("Codes Distribution "+h.__name__)
13     plt.show()
```

B Implémentation d'une table de hachage

On souhaite créer une table de hachage d'après un fichier qui recense les capitales des pays du monde entier. Cette table possède donc des clefs de type `str` (le pays) et des valeurs de type `str` (la capitale).

- B1. Écrire une fonction `import_csv()` qui importe les données du fichier `"capitals.csv"`. Cette fonction renvoie une liste de tuples (pays, capitale).

Solution :

```
1  def import_csv(filename):
2      # conventions :
3      # -- data are strings
4      with open(filename, "r") as f:
5          data = []
6          headers = f.readline().split(",")
7          for line in f:
8              words = line.split(',')
9              data.append((words[0], words[1])) # (country, capital)
10     return data
```

- B2. Écrire une fonction de prototype `init_hash_table(elements, table_size)` qui renvoie une table de hachage initialisée d'après le paramètre `elements`. Ce paramètre est la liste de tuples créée à la question précédente.

Solution :

```
1 def init_hash_table(elements, table_size):
2     hashtable = [[] for i in range(table_size)]
3     for key, value in elements:
4         hashed_value = hm(key)
5         index = hashed_value % table_size
6         hashtable[index].append((key, value))
7     return hashtable
```

- B3. Écrire une fonction de prototype `get_value(table, table_size, input_key)` qui permet d'accéder à l'élément de clef `input_key` de la table de hachage `table`. Par exemple, `get_value(ht, "Italy")`, n renvoie "Roma".

Solution :

```
1 def get_value(table, table_size, input_key):
2     hashed_value = hm(input_key)
3     index = hashed_value % table_size
4     element = table[index]
5     for key, value in element:
6         if key == input_key:
7             return value
8     return None
```

- B4. Créer l'ensemble de toutes les clefs de la table pour lesquelles il existe une valeur, puis parcourir la table à partir de cet ensemble. Les capitales apparaissent-elles dans un ordre quelconque?

Solution : L'ordre est quelconque puisque la table de hachage est répartie uniformément à partir de codes pseudo-aléatoires.

```
1 data = import_csv("capitals.csv")
2 print(data)
3 ht = init_hash_table(data, TABLE_SIZE)
4 print(ht)
5
6 print("Access to capital of Italy --> ", get_value(ht, TABLE_SIZE, "Italy"))
7
8 keys = []
9 for i, e in enumerate(ht):
10     if e:
11         keys.append(i)
12 print(keys, len(keys))
13
14 for k in keys:
```

```
15 print(k, ht[k])
```

R Naturellement, si par la suite vous avez besoin d'une table de hachage, il faut utiliser le type `dict` de Python et ne pas réinventer la poudre!