

Rapport de projet

Projet pour le cours Algorithme pour les data science.

SILVIU MANIU

M2 : Data Science

Introduction :

Dans le cadre de notre projet, en cours d'algorithmie nous avons étudiés plusieurs algorithmes et notre projet se concentre sur le min hashing ainsi que Locality Sensitive Hashing.

Le but est de pouvoir comparer des documents et avoir un taux de similarité entre plusieurs documents soumis en entrée.

Aujourd'hui plusieurs bibliothèques permettent d'avoir accès directement à ces algorithmes sans forcément en comprendre les fondements, mais nous avons vu en cours comment ces algorithmes cohabitent. Nous allons mettre tous cela en pratique lors de ce projet.

Le projet se compose en 2 parties : Une partie concernant le min hashing de base et une autre concernant le LSH.

Sommaire :

1. Cleaning
2. K-Shingles
3. Minhashing
4. LSH
5. Conclusion

1. Cleaning

Le projet nous donnait un fichier txt du quel on a extrait les données suivantes sous forme de DataSet ressemblant à ça a la base :

	Text_line
0	@stellargirl I looooooooooooooooooooo my Kindle2. ...
1	Reading my kindle2... Love it... Lee child's i...
2	Ok, first assesment of the #kindle2 ...it fuck...
3	@kenburbary You'll love your Kindle2. I've had...
4	@mikefish Fair enough. But i have the Kindle2...
...	...
492	Ask Programming: LaTeX or InDesign?: submitted...
493	On that note, I hate Word. I hate Pages. I hat...
494	Ahhh... back in a "real" text editing environm...
495	Trouble in Iran, I see. Hmm. Iran. Iran so far...
496	Reading the tweets coming out of Iran... The w...

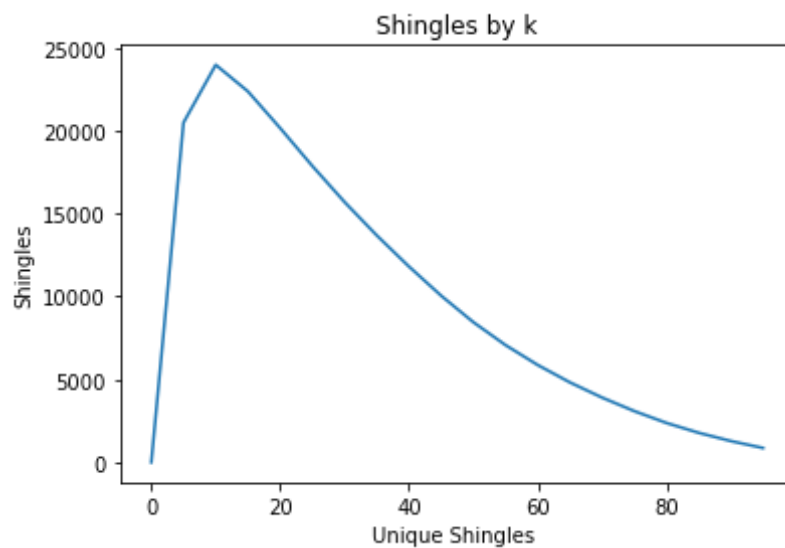
On commence au début par clean les données que l'on a en entrée car ces fichiers sont des tweets avec des répétitions, des emoji et afin de traiter le sens de ces informations avec une intelligence artificiel nous avons besoin uniquement du sens des tweets.

Ces données qui l'on veut clean deviennent sont documents chargés dans un tableau où les indices sont les identificateurs de documents et les valeurs le contenu des documents, c'est-à-dire le texte. Ensuite, on supprime les espaces entre les documents de sorte que chaque document devienne une longue chaîne de caractères sans espaces, c'est-à-dire une succession de caractères alphanumériques. Pour cela un dictionnaire est le format le plus adapté.

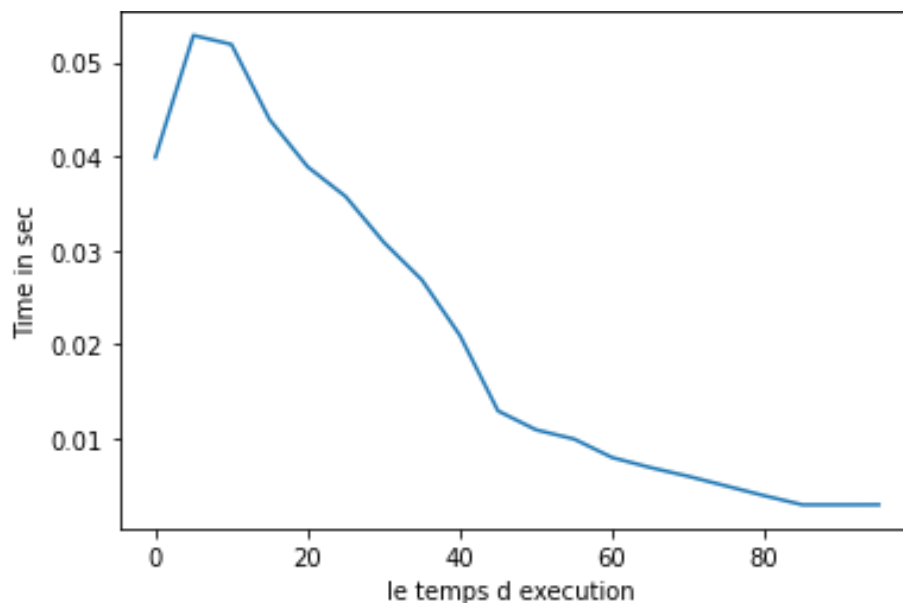
2. K-Shingles

Dans cette partie-là qui est une suite logique à la partie 1, nous divisons nos données en plusieurs fenêtre de mots faisant la même taille. Afin d'atteindre cet objectif, nous prenons comme référence les lab, ou la séparation en k était déjà imposé.

Pour cela une courbe affichant les shingles en fonction des entrées est ci-dessous.



En voulant visualiser son évolution, en fonction du temps d'exécution, nous avons décidé d'afficher le temps d'exécution de la fonction entre son démarrage et son return.



3. MinHashing

Pour construire la matrice de mini hachage, nous créons n fonctions de hachage. La création de ces fonctions de hachage se fait simplement par permutation des différents kshingles générés lors de l'étape précédente. Nous utiliserons ces fonctions de hachage pour construire la matrice de mini hachage qui sera alors de la taille `numéro_de_hash_function * nombre_de_documents`.

Lors de l'exécution du code voici ma matrice après un hachage, on peut voir que la différence avec l'affichage principal est notable :

```
[{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79}, {29, 30, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114}, {128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 29, 30, 80, 81, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127}, {28, 29, 30, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247}, {256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336}, {337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 122}, {380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000}
```

```
[{79, 108, 121, 207, 248, 311, 355, 394, 403, 427, 500, 568, 651, 662, 759, 858, 759, 982, 1028, 1036, 1102, 1124, 1175, 1189, 1301, 1355, 1398, 1448, 1429, 1552, 1668, 344, 1796, 1851, 2031, 2047, 2068, 2077, 2104, 2191, 2279, 2311, 2314, 2357, 2408, 2394, 2454, 1634, 816, 2536, 2563, 2608, 2600, 2653, 1189, 2719, 2823, 2857, 2867, 2902, 2934, 2934, 3032, 2934, 3068, 3091, 2934, 3198, 3206, 3310, 3345, 3369, 3414, 3460, 3488, 3574, 3596, 3683, 3723, 3738, 3902, 2222, 4069, 4125, 2125, 4196, 4202, 4236, 4268, 4358, 2031, 4441, 4536, 4581, 4196, 4708, 4712, 4778, 4817, 4868, 4915, 4954, 5032, 5126, 5132, 5205, 5262, 5309, 5347, 5377, 5377, 5426, 5478, 5513, 5548, 5574, 5625, 5647, 5666, 5713, 5740, 5796, 5899, 5986, 6017, 6029, 6039, 6084, 6039, 6029, 6199, 6212, 6241, 53, 2608, 6401, 6413, 6502, 6546, 6649, 6696, 6751, 6798, 6854, 6909, 6961, 6989, 867, 7106, 7126, 7191, 7215, 7273, 6909, 7340, 7350, 7452, 7537, 7574, 7637, 7678, 7716, 858, 7832, 7940, 7986, 8019, 8051, 7977, 8066, 8114, 8146, 8188, 8231, 8283, 8313, 8326, 8420, 8437, 8480, 3470, 8505, 8514, 8553, 8567, 8603, 8618, 655, 8676, 8686, 8648, 8710, 8789, 8877, 8894, 8947, 8983, 9059, 9151, 9155, 9194, 9229, 9342, 9419, 9459, 9508, 9625, 9678, 277, 6961, 9808, 9851, 867, 9978, 10028, 10066, 10096, 10145, 10261, 2608, 10348, 10414, 10450, 10451, 10477, 10507, 10517, 10545, 10553, 10587, 10631, 5625, 3196, 10775, 6033, 10861, 10875, 28, 10922, 10976, 10990, 11028, 11058, 11086, 11086, 6273, 11140, 11226, 11242, 11359, 11378, 11440, 965, 11496, 11512, 11541, 11614, 11646, 4498, 11687, 11737, 11804, 11822, 11868, 273, 11964, 12011, 12067, 10555, 12108, 12162, 12178, 12215, 12245, 12280, 12295, 12385, 166, 12462, 12513, 8223, 12557, 12586, 2431, 12557, 12647, 12655, 2637, 12697, 12710, 12722, 12769, 12850, 4915, 2719, 79, 12949, 12245, 13010, 13028, 3029, 13057, 13110, 918, 13110, 13158, 13200, 13266, 858, 13345, 13368, 13403, 7340, 8909, 13573, 4565, 13657, 13666, 13074, 13723, 395}
```

Lors de cette partie-là, le Lab1 a été d'une grande aide car celui-ci donne les bases essentielles à la compréhension du Minhash et des shingles

4. LSH

Dans cette partie, nous avons mis en œuvre la fonction de hachage sensible à la localité LSH.

Le principe de la LSH est le suivant : diviser notre matrice signature en différentes bandes b contenant plusieurs r lignes chacune. On apprend après que $br = n$, la taille des signatures. A l'intérieur d'une bande particulière, on se retrouve avec L « sous signature » de taille r associés à L codes génétiques. Dans chacune des b bandes on va « hasher » ces sous signatures vers un nombre relativement restreint de buckets (on choisit généralement un million de buckets) .

Pour le hachage sensible aux localités, les performances varient en fonction sur le seuil de similarité et le nombre de bandes que nous définissons.

Pour cela nous avons défini 2 fonctions : `precalcul` & `bandadaptor` qui sont là pour séparer la matrice en différentes bande et calculer la hashage à l'intérieur. Mon but au début était d'implémenter le Minhash mais vu plusieurs problème . J'ai décider de hasher d'une autre manière.

Avec LSH, nous avons une bibliothèque qui permet d'appliquer automatiquement ces fonctions-là, du coup en important directement `Textreuse` on a accès aux fonctions `Lsh` , `Minhash` .. Nos calculs sont beaucoup plus simples comme ça, voici un exemple de cette implémentation.

```
Entrée [ ]: import textreuse

n_param=100
seed_p = 330
b_param=20

buckets = lsh(x = S,
bands = b_param,
progress = TRUE)
candidates_pairs = lsh_candidates(buckets = buckets)

similarity_candidate_pairs = lsh_compare(candidates=candidates_pairs,
corpus=S_matrix,
f=jaccard_similarity,
progress = TRUE)

threshold=0.9
similar_pairs=similarity_candidate_pairs[similarity_candidate_pairs$score>threshold,]
```