

Đề thi tuyển sinh Trại hè Toán học và Ứng dụng 2021 —PHẦN TỰ NGHIÊN CỨU—

Cấu trúc đề thi và hướng dẫn

Đây là phần Đề thi tuyển sinh phần Câu hỏi tự nghiên cứu cho trại hè PiMA 2021. Đề thi gồm nhiều phần câu hỏi xoay quanh chủ đề “Phân cụm K-means” trong Khoa học Dữ liệu, và đòi hỏi từ bạn kỹ năng tự tìm kiếm thông tin và tiếp thu kiến thức. Bạn có thể trình bày câu trả lời của mình bằng nhiều hình thức tùy ý, nhưng hãy đảm bảo đủ ý và không quá dài dòng. Một cách tìm kiếm hiệu quả là gõ tên các khái niệm vào thanh tìm kiếm của trang web www.google.com. Ngoài ra, phần phụ lục cuối đề thi sẽ hỗ trợ bạn về nền tảng để lập trình Python.

Thuật ngữ tiếng Anh

Một số thuật ngữ chúng mình sẽ dùng từ tiếng Anh thay cho tiếng Việt ở một hoặc nhiều chỗ trong file đề bài này:

STT	Tiếng Việt	Tiếng Anh	STT	Tiếng Việt	Tiếng Anh
1	Phân cụm	Clustering	2	Cụm	Cluster
3	Tập dữ liệu	Dataset	4	Điểm dữ liệu	Data point
5	Trọng tâm	Centroid	6	Nhãn	Label
7	Đầu vào	Input	8	Đầu ra	Output
9	Cho ra (kết quả)	Return	10	Chú giải hoá (code)	Comment
11	Ngừng chú giải hoá	Uncomment			

Phần câu hỏi chính - Phân cụm K-means

Thuật toán phân cụm **K-means** là một thuật toán phổ biến đối với người nhập môn Khoa học Dữ liệu. Bạn hãy tìm hiểu thật nhiều về thuật toán này để trả lời các câu hỏi sau.

Câu hỏi lý thuyết

Trả lời các câu hỏi sau và giải thích ngắn gọn câu trả lời của bạn.

LT1 Mô tả bài toán mà thuật toán này muốn giải. **Gợi ý:** trả lời các câu hỏi nhỏ:

- Dữ liệu đầu vào (**input**), đầu ra (**output**) của thuật toán này là gì?
- Thuật toán này muốn tối đa / tối thiểu hoá đại lượng gì?

LT2 Mô tả thuật toán K-means bằng ngôn ngữ của bạn. **Gợi ý:** trả lời các câu hỏi nhỏ:

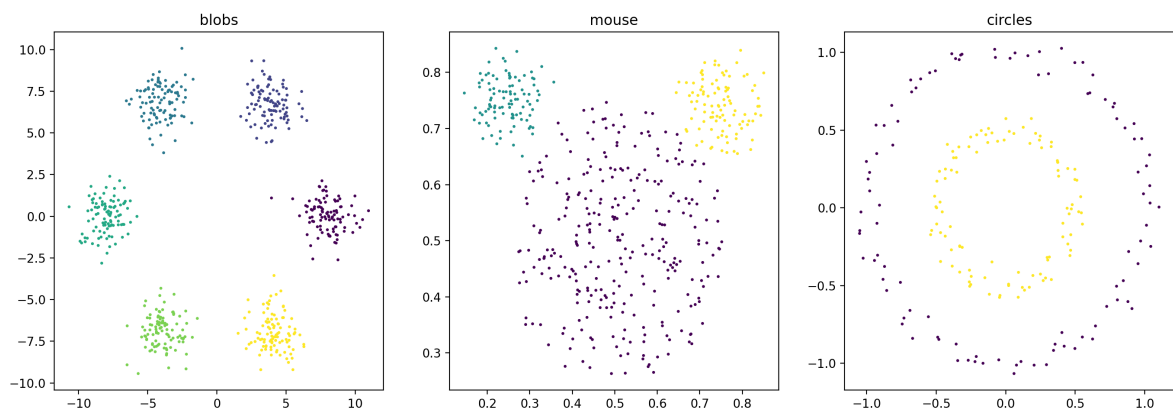
- Viết ra các bước của thuật toán theo trình tự, dưới dạng danh sách hoặc các sơ đồ như sơ đồ khối.
- Giải thích mục đích từng bước trong thuật toán một cách ngắn gọn.

Mở trang web sau: [K-means visualization by N. Harris](#). Bạn có thể nhấn vào từng nút bấm trên trang web để hiểu cách K-means được thể hiện bằng đồ hoạ. Sau khi đã hiểu cách trang web hoạt động, trả lời các câu hỏi sau:

LT3 Đối với tập dữ liệu (**dataset**) “Gaussian Mixture” và số cụm là 3, tìm 1 cách chọn 3 trọng tâm ban đầu sao cho output của thuật toán không phải là tối ưu. Giải thích tại sao lại có hiện tượng này.

LT4 Hãy thử đưa ra một thuật toán, gọi là “Centroids Smart Select”, để chọn các trọng tâm ban đầu sao cho tăng khả năng tìm chính xác các cụm. Mô tả input, output và các trình tự các bước của thuật toán đó. **Gợi ý:** Bạn có thể tìm thấy vài thuật toán như vậy trên mạng.

LT5 Các biểu đồ sau thể hiện một số datasets gồm các điểm trong mặt phẳng \mathbb{R}^2 . Các điểm cùng một cụm được tô cùng màu. Với mỗi datasets, dự đoán xem thuật toán K-means, với số cụm đúng được cho trước, có thể tìm chính xác các cụm được không. Giải thích tại sao. **Gợi ý:** Hãy thử thuật toán này với các dataset còn lại ngoài “Gaussian Mixture” trên trang web của N. Harris và quan sát kết quả.



Câu hỏi lập trình

Lưu ý: Nếu bạn chưa có kinh nghiệm lập trình, có thể làm phần này cuối cùng hoặc bỏ qua nếu thời gian không cho phép. Có rất nhiều trại sinh PiMA chưa có kinh nghiệm lập trình tại thời điểm tham gia trại.

Trong folder này, bạn sẽ tìm thấy các file sau:

1. File đề bài này.
2. `my_kmeans.py` - chứa một đoạn code bằng Python triển khai thuật toán K-means. Trong các câu hỏi TH1 và TH2 bạn sẽ phải hoàn thành file này.
3. `main.py` - chứa một đoạn code bằng Python để kiểm tra đoạn code bạn viết ở file `my_kmeans.py`. Chạy file này khi bạn muốn kiểm tra code của mình.
4. `blobs.csv` - dataset "blobs", `mouse.csv` - dataset "mouse", `circles.csv` - dataset "circles", `emails.csv` - dataset "emails", `senate.csv` - dataset "lịch sử bỏ phiếu của Thượng viện Mỹ khóa 114".

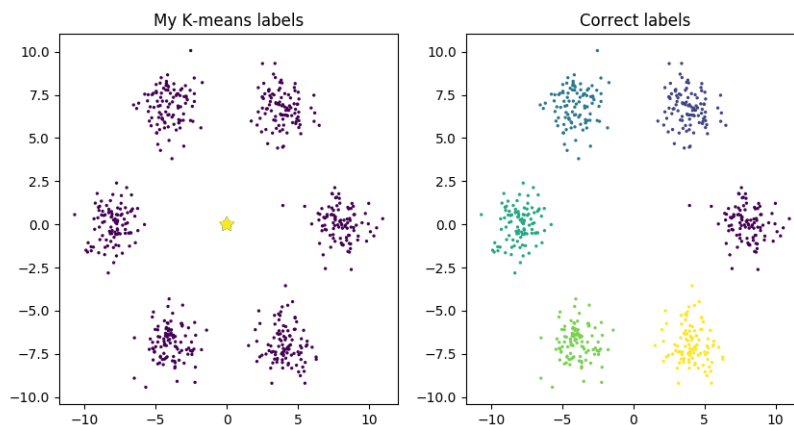
Hãy **uncomment** dòng số 72 ở `main.py` bằng cách xóa dấu `#` ở đầu dòng:

```
1 # kmeans_test('blobs', smart=False) # Uncomment de test cau hoi TH1
```

thành

```
1 kmeans_test('blobs', smart=False) # Uncomment de test cau hoi TH1
```

Chạy file `main.py`, bạn sẽ thấy kết quả như sau:



Biểu đồ bên trái thể hiện output của thuật toán K-means trong file `my_kmeans.py` với các điểm data cùng cluster được tô cùng màu. Biểu tượng ngôi sao thể hiện vị trí trọng tâm (centroid) của các cluster.

Biểu đồ bên phải thể hiện các điểm data với màu tương ứng với cluster chính xác.

Vì thuật toán K-means **chưa được triển khai đúng**, cụ thể là các hàm `find_nearest_centroid`, `find_new_centroids` (và `centroids_smart_select` cho phiên bản cải tiến) trong file `my_kmeans.py`, output chỉ có một cluster, so với 6 cluster chính xác.

Để **comment** dòng này lại, chỉ cần thêm dấu `#` lại vào đầu dòng. Khi chạy file, các dòng được comment (tức là có dấu `#` đầu dòng) sẽ không được tính là một phần của code.

Lưu ý: Chỉ xóa dấu `#` đầu dòng, hãy giữ nguyên dấu `#` thứ hai.

Qua các câu hỏi sau, bạn sẽ điền code của mình vào để hoàn thiện thuật toán K-means.

Cảnh báo: Bạn không được phép thay đổi bất kỳ chỗ nào trong các file, trừ ba hàm `find_nearest_centroid`, `find_new_centroids` và `centroids_smart_select` trong `my_kmeans.py`

TH1 Hoàn thành đoạn code triển khai K-means bằng cách điền nội dung vào các hàm `find_nearest_centroid` và `find_new_centroids` trong `my_kmeans.py`.

Uncomment dòng `kmeans_test('blobs', smart=False)` và chạy file `main.py` để kiểm tra lỗi và thử output của đoạn code K-means của bạn.

Sau khi hoàn thành câu TH1 hãy comment dòng này lại.

Uncomment dòng `kmeans_testblobs(100, smart=False)` và chạy `main.py`. Dòng này sẽ áp dụng thuật toán K-means mà bạn triển khai lên dataset “blobs” 100 lần, mỗi lần tính một số điểm (score) bằng số điểm dữ liệu được phân đúng cụm. Dòng cuối sẽ hiển thị score trung bình của 100 lần chạy.

TH2 Hoàn thành việc triển khai thuật toán “Centroid Smart Select” mà bạn nghĩ ra trong câu hỏi LT4 bằng cách điền nội dung vào hàm `centroids_smart_select` trong file `my_kmeans.py`. Đổi dòng `kmeans_testblobs(100, smart=False)` thành `kmeans_testblobs(100, smart=True)` và chạy `main.py` để thử output của thuật toán `centroids_smart_select` của bạn. So sánh output khi `smart=True` so với khi `smart=False` và giải thích nguyên nhân cho sự khác nhau.

Sau khi hoàn thành TH2 hãy comment dòng này lại.

TH3 Uncomment dòng `kmeans_test('mouse', smart=False)`, chạy file `main.py` và mô tả output. Đổi dòng này thành `kmeans_test('mouse', smart=True)`, chạy và so sánh output với khi `smart=False`. Lặp lại 2 bước trên với dataset `'circles'`. Thuật toán `centroid_smart_select` có cải thiện output của K-means trên các dataset “mouse” và “circles” không?

Câu hỏi nâng cao

Với mỗi câu hỏi sau, bạn sẽ được hỏi liệu thuật toán phân cụm K-means có thể chia dataset tương ứng thành các cụm theo yêu cầu không. Giải thích ngắn gọn liệu bạn nghĩ K-means có thể được áp dụng vào mỗi trường hợp. Nếu không áp dụng được, hãy đề xuất một biến thể / mở rộng của K-means mà bạn nghĩ có thể áp dụng được cho bài toán. Bạn có thể nộp code của mình nếu có (**không bắt buộc**). Ở phần này, bạn được phép dùng các phiên bản K-means có trong các thư viện online, ví dụ như `sklearn.cluster.KMeans`.

NC4 File `senate.csv`¹ chứa lịch sử bỏ phiếu của các thượng nghị sĩ Mỹ khoá 114 đối với các dự luật. Ví dụ 3 dòng đầu của dataset:

¹Bạn cũng có thể download file này tại [đây](#)



```
1 name,party,state,00001,00004,00005,00006,00007,00008,00009,00010,00011,00012,00013,00014,00015,00016,00017,00018,00019,00020,00021,00022,00023,00024,00025,00026,00027,00028,00029,00030,00031,00032,00033,00034,00035,00036,00037,00038,00039,00040,00041,00042,00043,00044,00045,00046,00047,00048,00049,00050,00051,00052,00053,00054,00055,00056,00057,00058,00059,00060,00061,00062,00063,00064,00065,00066,00067,00068,00069,00070,00071,00072,00073,00074,00075,00076,00077,00078,00079,00080,00081,00082,00083,00084,00085,00086,00087,00088,00089,00090,00091,00092,00093,00094,00095,00096,00097,00098,00099,00100,00101,00102,00103,00104,00105,00106,00107,00108,00109,00110,00111,00112,00113,00114,00115,00116,00117,00118,00119,00120,00121,00122,00123,00124,00125,00126,00127,00128,00129,00130,00131,00132,00133,00134,00135,00136,00137,00138,00139,00140,00141,00142,00143,00144,00145,00146,00147,00148,00149,00150,00151,00152,00153,00154,00155,00156,00157,00158,00159,00160,00161,00162,00163,00164,00165,00166,00167,00168,00169,00170,00171,00172,00173,00174,00175,00176,00177,00178,00179,00180,00181,00182,00183,00184,00185,00186,00187,00188,00189,00190,00191,00192,00193,00194,00195,00196,00197,00198,00199,00200,00201,00202,00203,00204,00205,00206,00207,00208,00209,00210,00211,00212,00213,00214,00215,00216,00217,00218,00219,00220,00221,00222,00223,00224,00225,00226,00227,00228,00229,00230,00231,00232,00233,00234,00235,00236,00237,00238,00239,00240,00241,00242,00243,00244,00245,00246,00247,00248,00249,00250,00251,00252,00253,00254,00255,00256,00257,00258,00259,00260,00261,00262,00263,00264,00265,00266,00267,00268,00269,00270,00271,00272,00273,00274,00275,00276,00277,00278,00279,00280,00281,00282,00283,00284,00285,00286,00287,00288,00289,00290,00291,00292,00293,00294,00295,00296,00297,00298,00299,00300,00301,00302,00303,00304,00305,00306,00307,00308,00309,00310,00311,00312,00313,00314,00315,00316,00317,00318,00319,00320,00321,00322,00323,00324,00325,00326,00327,00328,00329,00330,00331,00332,00333,00334,00335,00336,00337,00338,00339,00340,00341,00342,00343,00344,00345,00346,00347,00348,00349,00350,00351,00352,00353,00354,00355,00356,00357,00358,00359,00360,00361,00362,00363,00364,00365,00366,00367,00368,00369,00370,00371,00372,00373,00374,00375,00376,00377,00378,00379,00380,00381,00382,00383,00384,00385,00386,00387,00388,00389,00390,00391,00392,00393,00394,00395,00396,00397,00398,00399,00400,00401,00402,00403,00404,00405,00406,00407,00408,00409,00410,00411,00412,00413,00414,00415,00416,00417,00418,00419,00420,00421,00422,00423,00424,00425,00426,00427,00428,00429,00430,00431,00432,00433,00434,00435,00436,00437,00438,00439,00440,00441,00442,00443,00444,00445,00446,00447,00448,00449,00450,00451,00452,00453,00454,00455,00456,00457,00458,00459,00460,00461,00462,00463,00464,00465,00466,00467,00468,00469,00470,00471,00472,00473,00474,00475,00476,00477,00478,00479,00480,00481,00482,00483,00484,00485,00486,00487,00488,00489,00490,00491,00492,00493,00494,00495,00496,00497,00498,00499,00500,00501,00502,00503,00504,00505,00506,00507,00508,00509,00510,00511,00512,00513,00514,00515,00516,00517,00518,00519,00520,00521,00522,00523,00524,00525,00526,00527,00528,00529,00530,00531,00532,00533,00534,00535,00536,00537,00538,00539,00540,00541,00542,00543,00544,00545,00546,00547,00548,00549,00550,00551,00552,00553,00554,00555,00556,00557,00558,00559,00560,00561,00562,00563,00564,00565,00566,00567,00568,00569,00570,00571,00572,00573,00574,00575,00576,00577,00578,00579,00580,00581,00582,00583,00584,00585,00586,00587,00588,00589,00590,00591,00592,00593,00594,00595,00596,00597,00598,00599,00600,00601,00602,00603,00604,00605,00606,00607,00608,00609,00610,00611,00612,00613,00614,00615,00616,00617,00618,00619,00620,00621,00622,00623,00624,00625,00626,00627,00628,00629,00630,00631,00632,00633,00634,00635,00636,00637,00638,00639,00640,00641,00642,00643,00644,00645,00646,00647,00648,00649,00650,00651,00652,00653,00654,00655,00656,00657,00658,00659,00660,00661,00662,00663,00664,00665,00666,00667,00668,00669,00670,00671,00672,00673,00674,00675,00676,00677,00678,00679,00680,00681,00682,00683,00
```

Dòng đầu chứa tên các biến `name` (tên), `party` (đảng), `state` (bang) và lịch sử bỏ phiếu cho các dự luật có mã từ `00001` đến `00047`. Từ dòng thứ hai trở đi, mỗi dòng chứa tên, đảng, bang và lịch sử bỏ phiếu của từng thượng nghị sĩ. Số `1` là Thuận, `0` là Chống và `0.5` là phiếu Trắng. Ví dụ, dòng thứ ba cho biết: Bà Ayotte là thượng nghị sĩ đảng Republican (R) từ bang New Hampshire (NH), bà bỏ phiếu Thuận cho các dự luật có mã 4, 5, 6, 7, 10, 26, 38, 44 và bỏ phiếu Chống cho các dự luật còn lại.

Câu hỏi: Liệu có thể dùng K-means để chia các thượng nghị sĩ thành 2 cụm sao cho các thành viên cùng đảng vào chung cụm nhiều nhất có thể? Hơn nữa, liệu có thể dùng K-means để tìm ra các phe phái nhỏ trong từng đảng với lịch sử bỏ phiếu tương tự nhau? Nếu không, bạn có thể đề xuất một cải biến của K-means để giải hai bài toán trên không?

NC5 Công ty PiMA Corp gồm 100 nhân viên được chia làm 4 phòng ban: Nghiên cứu, Kỹ sư, Marketing và Nhân sự. Mỗi nhân viên có xu hướng trao đổi nhiều emails với người cùng phòng mình hơn người khác phòng. File `emails.csv` chứa một ma trận $M \in \mathbb{R}^{100 \times 100}$, trong đó M_{ij} là số emails nhân viên i gửi tới nhân viên j trong vòng 1 tháng. **Lưu ý:** M_{ij} có thể khác M_{ji} .

Câu hỏi: Liệu có thể dùng K-means để tìm ra ai thuộc phòng ban nào? Nếu không, bạn có thể đề xuất một cải biến của K-means để giải bài toán trên không?

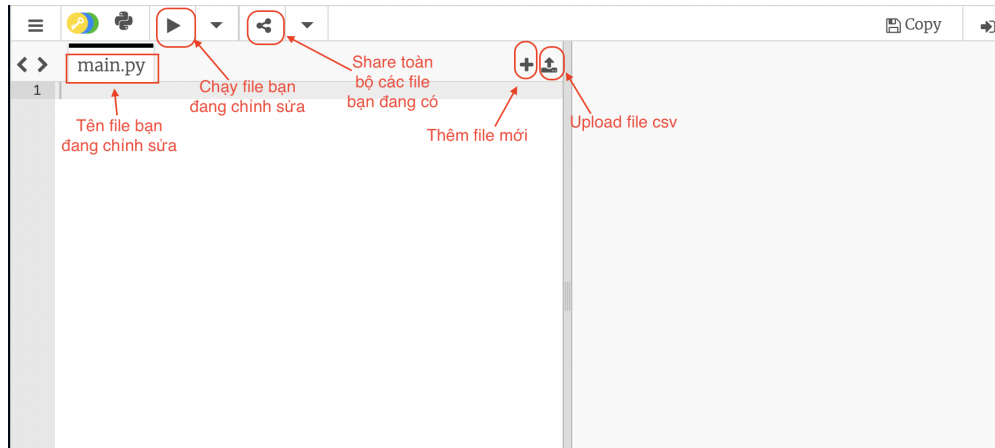
A. Hướng dẫn thêm về phần Lập trình

Bạn có thể tìm hiểu thêm về Python tại [Laptoprinh.io](https://laptoprinh.io), [Tutorialspoint](https://tutorialspoint.com), hoặc sử dụng Google.

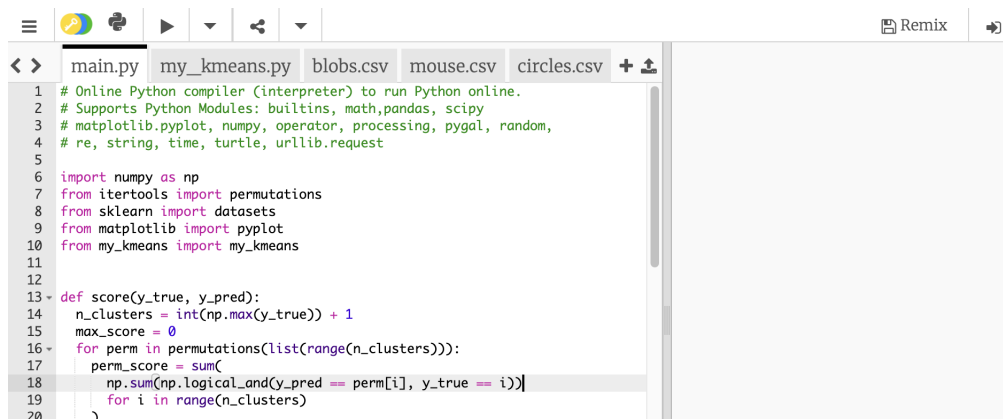
Nếu bạn không có Python trong máy tính, hoặc bản Python của bạn không có các thư viện `numpy`, `scipy`, `sklearn`, v.v. thì có thể dùng trang web [Trinket.io](https://trinket.io). Bạn có thể dùng bất cứ nền tảng chạy code Python nào tùy thích.

Cách sử dụng Trinket.io

Vào link [Trinket.io](https://trinket.io), bạn sẽ thấy các nút thao tác sau:



Copy và paste toàn bộ code trong file `main.py` vào, sau đó tạo thêm một file mới với tên `my_kmeans.py`. Copy và paste toàn bộ code trong file `my_kmeans.py` trong folder này vào file vừa tạo. Upload các file `blobs.csv`, `mouse.csv`, `circles.csv` bằng nút upload. Bạn sẽ thấy như sau:



Bấm vào nút chạy code (hình ►) để quan sát kết quả.

—HẾT—