

Proyecto 3

1st Barrios, Alonso

Universidad de Tecnología e Ingeniería
Ciencia de la Computación
Lima, Perú
alonso.barrios@utec.edu.pe

2nd Rodriguez, Mauricio

Universidad de Tecnología e Ingeniería
Ciencia de la Computación
Lima, Perú
mauricio.rodriguez@utec.edu.pe

3rd Tenazoa, Renzo

Universidad de Tecnología e Ingeniería
Ciencia de la Computación
Lima, Perú
renzo.tenazoa@utec.edu.pe

Abstract—Este proyecto está desarrollado para poner a prueba lo aprendido acerca de los modelos de Clustering. El objetivo del proyecto es agrupar distintos puntos de un dataset, para predecir su clasificación haciendo uso de distintos modelos de clustering y analizar las graficas de error para definir cual es el que mejor resultado aporta.

Index Terms—clustering, predicción, error, underfitting, overfitting

I. INTRODUCCIÓN

Con el objetivo de poner a prueba los modelos de *clustering* se utilizó el dataset *datasettissue.txt*, en cual consta de 22215 filas y 189 columnas que representan las distintas características de diferentes tipos de tejidos. Se buscó predecir el tipo de tejido al cual pertenece cada uno de los puntos, mediante agrupación haciendo uso de modelos de *clustering* tales como *GMM*, *KMeans*, *DBScan* y *Agglomerative Hierarchical Clustering*. Posteriormente, se comparó el resultado obtenido de los modelos con los presentados en el dataset para así lograr presentar las matrices de confusión y analizar cual de los modelos de *clustering* nos brinda mejores resultados.

II. EXPLICACIÓN

A. KMeans

Para este método se usaron las siguientes funciones:

- *distance()*: Obtiene la distancia euclidiana entre dos elementos del dataset.
- *initialize_centroids()*: Escoge elementos del dataset de forma aleatoria y los utiliza para inicializar los centroides.
- *_create_clusters()*: Asocia en clusters todos los índices de los datos que tengan el mismo centroide más cercano.
- *_get_closest_centroid(row)*: función de ayuda utilizada en *create_clusters()*. Dado un elemento del dataset devuelve el centroide más cercano.
- *_update_centroids()*: Se encarga de calcular el error de cada fase de entrenamiento.
- *_update_centroids()*: Actualiza el centroide a la media de todos los elementos del cluster relacionados a este centroide.
- *kmeans()*: Es la función principal, inicializa los centroides con la función *initialize_centroids*. Luego de ello, genera los clusters asociados a los centroides actualiza estos últimos para ubicarse en la media del cluster. Este proceso se da con las funciones *_create_clusters()*, *_update_centroids()*, y se realiza hasta que se llegue

un numero determinado de iteraciones o hasta que los centroides convergan. Finalmente retorna los clusters.

B. Sklearn

Para el resto de métodos de *clustering* se usó la librería *sklearn*.

III. FUNCIONES GENERALES

- *reduction()*: Se encarga de reducir la dimensionalidad del dataset a la dimensión deseada.
- *normalize()*: Se encarga de normalizar los datos.
- *toCSV()*: Se encarga de pasar los datos de los archivos txt a formato csv para hacer uso eficiente de la lectura por disco.
- *read()*: Se encarga de los archivo csv proporcionados por la función *toCSV*.

IV. EXPERIMENTOS

Para la parte de experimentación se contó con datos de 22215 atributos, por lo que se usó un método de reducción de dimensión llamado *PCA* (*principal component analysis*). Luego de diversas pruebas se logró reducir la cantidad de atributos a 20 sin perder demasiado el cubrimiento total de los datos. Posteriormente, se comparó los *labels* obtenidos por cada modelo de *clustering* y se comparó con los del dataset para así analizar las gráficas de errores de cada modelo.

A. GMM

En este modelo se usaron los siguientes parámetros:

- *n_components*: Como 7, ya que este parámetro nos indica la cantidad de *clusters* que queremos tener. Ya que la cantidad de nombres de los tejidos es 7, creamos la misma cantidad de *clusters*.
- *random_state*: Como 0, este parámetro indica la semilla que queremos que tenga el modelo, como se desea replicar los resultados en diferentes computadoras se colocó un número estándar.

Los resultados obtenidos fueron:

Fig. 1. GMM

B. KMeans

- *data*: Este parámetro simplemente recibe los elementos "X" de la *data*. En este caso recibe la data normalizada y posteriormente reducida mediante análisis de componentes principales(PCA)
- *k*: Como 7, ya que este parámetro nos indica la cantidad de *clusters* que queremos tener. Ya que la cantidad de nombres de los tejidos es 7, creamos la misma cantidad de *clusters*.
- *iterations*: Como 100 de forma estándar, este parámetro determina el número máximo de iteraciones que realizará el algoritmo. Es decir el número de actualizaciones que tendrán los centroides, y por ende, los clusters.

[illegible]

Fig. 2. Impresión de clusters obtenidos

- *eps*: Como 3, este parámetro nos indica el radio que queremos usar para capturar a los vecinos más cercanos.
- *min_samples*: Como 2, este parámetro indica el mínimo número de elementos que tienen que tener los *clusters*.

```

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:
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:
115:
116:
117:
118:
119:
120:
121:
122:
123:
124:
125:
126:
127:
128:
129:
130:
131:
132:
133:
134:
135:
136:
137:
138:
139:
140:
141:
142:
143:
144:
145:
146:
147:
148:
149:
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:
248:
249:
250:
251:
252:
253:
254:
255:
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:
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:
8
```

Fig. 3. DBScan

CONCLUSIONES

- Analizando los distintos resultados que nos otorga cada modelo de *clustering*, el que nos da mayor *accuracy* es el modelo de
- Depende de los hiperparámetros que le pasemos a cada modelo para aumentar su precisión de clasificación.
- Es necesario la reducción de datos ya que si no los modelos serían muy pesados computacionalmente.
- Al no usar todos los atributos es que no se genera un mayor *accuracy* en los modelos.
- Al disminuir el *alpha* el algoritmo necesitará mayor cantidad de *epoch* para poder hallar un error mínimo óptimo.
- Para obtener un buen resultado con el modelo DBScan es un prueba y error de los hiperparámetros que le podríamos pasar.

REPOSITORIO

Enlace al repositorio de Github:
<https://github.com/mauricio-rodriguez/ai-project-3>

REFERENCES

- [1] Montero, R. (2016). *Modelo de regresión lineal múltiple*. Universidad de Granada. España.
- [2] Bishop, C. (2006). *Pattern Recognition and Machine Learning*. Cambridge CB3 0FB, U.K.
- [3] Jain, S. (2020). *Some Notes on DBSCAN Algorithm*. Recovery from: <https://medium.com/analytics-vidhya/some-notes-on-dbscan-algorithm-61a2e9acce29>

resultados es que la mayoría de datos nos pertenecen a ningún