

Lab01: k-Nearest Neighbors

Por:
Ian Gabriel Cañas Fernández, 1092228

Profesor: Juan S. Pérez R.,
Asignatura: INL367L, Secc 01

Resumen:

En este laboratorio estaremos implementando el algoritmo de k Nearest Neighbors (kNN), algoritmo que consiste en la clasificación automática de varios puestos de muestra respecto a los puntos conocidos más cercanos, para ello se ha iniciado haciendo un proceso analítico en el que se simulan los cálculos llevados a cabo y se concluye con la prueba de este mediante un set de 70 000 datos de entrenamiento y siendo probado con 20 000 datos de prueba.

Ejercicios previos:

P1.1 Predicción analítica de género.

En la Ilustración 1, Ilustración 2, Ilustración 3 y Ilustración 4, se presenta una síntesis de lo que se lleva a cabo en el algoritmo y los resultados obtenidos, en tal caso se ha optado por utilizar una distancia euclídeana, pues estamos asumiendo un peso equitativo entre todos los datos.

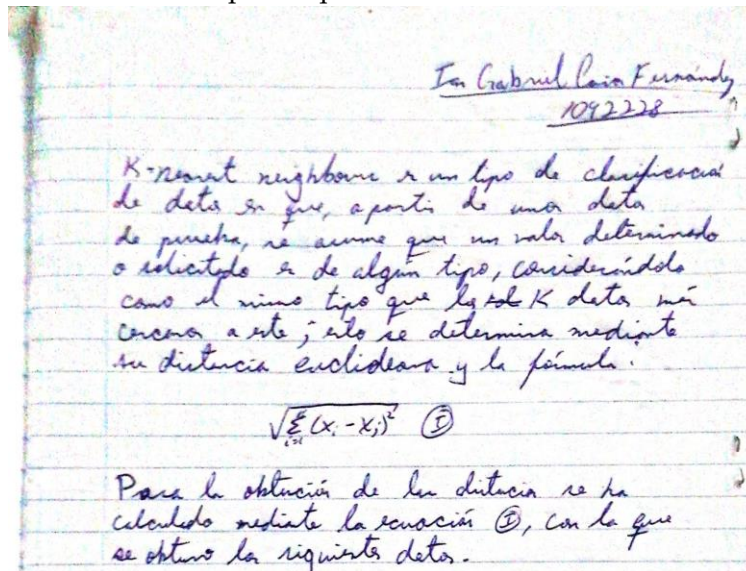


Ilustración 1. Proceso analítico llevado a cabo.

Handwritten table titled "tabla de Prueba" showing test data. The table has four columns: height_sample, weight_sample, alco_sample, and gender. There are 10 rows of data, numbered 1 to 10. Some cells contain handwritten marks like circles or checkmarks.

Test Set:	height_sample	weight_sample	alco_sample	gender
1	163	0	66	1
2	161	0	73	0
3	177	-	86	1
4	174	0	71	0
5	171	-	108	0
6	173	-	66	0
7	170	0	80	0
8	175	0	76	0
9	184	-	111	0
10	147	0	80	1

Ilustración 2. Tabla de datos de prueba.

Train Set:																
id	gender	height	weight	alco	Distancia 1	Distancia 2	Distancia 3	Distancia 4	Distancia 5	Distancia 6	Distancia 7	Distancia 8	Distancia 9	Distancia 10		
0	2	168	62	0	6.4	13.0	25.6	2	10.5	46.1	2	5.4	18.1	15.7	51.5	27.7
1	1	156	85	0	20.2	18.0	21.0	22.8	27.5	25.5	14.9	21.0	28.2	10.1	24.1	22.1
2	1	165	64	0	1	2.8	9.8	25.1	11.4	44.4	1	3.2	16.8	15.6	50.7	24.1
3	2	169	82	0	17.1	12.0	2	8.9	12.1	26.1	16.5	2	2	6.5	32.6	22.1
4	1	156	56	0	12.2	12.2	36.6	21.4	54.1	19.7	37.8	22.4	61.7	75.6	75.6	75.6
8	1	151	67	0	12.0	11.7	32.2	23.3	45.6	22.0	23.0	25.6	55.0	1	11.6	11.6
9	1	157	93	0	27.7	20.4	21.2	27.8	20.5	31.4	18.4	24.8	3	32.4	16.4	16.4
12	2	176	95	0	32.5	27.3	2	9.1	24.3	23.4	17.0	19.2	2	27.1	24.9	24.9
13	1	158	71	0	7.1	3.6	24.2	16.0	39.2	15.8	15.0	17.7	47.7	14.2	14.2	14.2
14	1	164	68	0	1	2.2	5.8	22.2	1	10.4	40.6	9.2	13.4	15.6	47.4	20.8
15	1	169	80	0	15.2	10.6	10.0	1	10.8	28.1	14.6	1	10.1	7.2	34.8	22.0
16	2	173	60	0	11.7	17.7	26.3	11.0	48.0	2	4.0	20.2	16.1	52.2	32.8	32.8
18	2	165	60	0	2	6.8	13.6	28.6	14.2	48.4	10.0	20.6	18.9	54.4	26.9	26.9
21	1	158	78	0	13.0	5.8	20.6	17.5	32.7	19.2	1	12.2	17.1	42.0	1	11.2
23	2	181	95	1	34.1	29.7	2	8.8	25.0	2	16.4	30.1	18.6	19.9	2	16.3
					1	1	2	1	2	2	1	1	2	1		

Ilustración 3. Comparación con distancias menores.

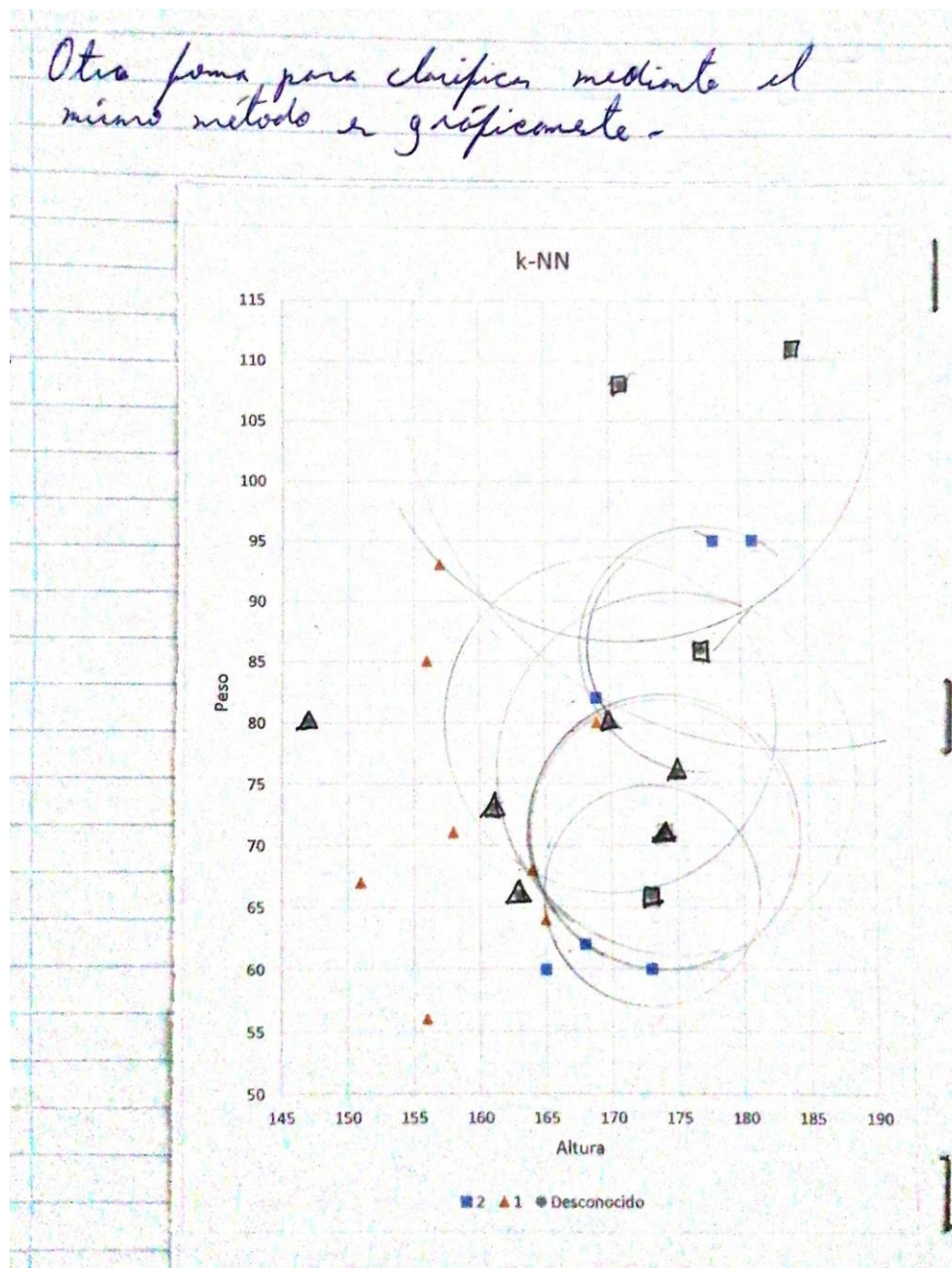


Ilustración 4. Resultados analíticos.

E1.1 Implementación en Matlab de kNN

Para el proceso computarizado de la implementación de kNN, se ha empezado importando los datos en dos tablas, luego, se ha identificado los datos para próximamente ser graficados y habiéndose solicitado la cantidad de vecinos con los que se solicita trabajar, véase la Ilustración 5.

```
1 clear, clc
2
3 train_set = readtable("Train set.xlsx"); %readtable("Train2.xlsx");
4 test_set = readtable("Test set.xlsx"); %readtable("Test2.xlsx");
5
6 g1 = train_set(train_set.gender==1, :);
7 g2 = train_set(train_set.gender==2, :);
8
9 plot(g1.height, g1.weight, '.b')
10 hold on
11 plot(g2.height, g2.weight, '.y')
12
13 k = input('Introduzca la cantidad de vecinos: ');
```

Ilustración 5. Clasificación de datos.

A continuación, se obtuvo la distancia de los puntos de prueba respecto al modelo, desde donde se comparó los k valores más pequeños y se obtuvo sus respectivos índices, mediante los que se podrá acceder y comparar con los resultados obtenidos, en la Ilustración 6 se puede observar dicho extracto de código.

```
15 nTrainData = length(train_set.height) %size(train_set);
16
17 for i = 1:length(test_set.height)
18     %i
19
20     Rep = repmat([test_set.height(i), test_set.weight(i)], nTrainData, 1);
21     d = ((Rep - [train_set.height(:) train_set.weight(:)])).^2);
22     d = sqrt(d(:,1)+d(:,2));
23
24     [dis pos] = sort(d,'ascend');
25     kNN=pos(1:k);
26     kND=dis(1:k);
27
```

Ilustración 6. Obtención de distancias.

Finalmente se analizan los datos obtenidos por saber la etiqueta que le tocaría a cada valor y finalmente se grafican, este extracto de código se presenta en la Ilustración 7:

```
29 c1 = 0;
30 c2 = 0;
31 for m = 1:k
32     if (train_set.gender(pos(m))==1)
33         c1 = c1+1;
34     elseif (train_set.gender(pos(m))==2)
35         c2 = c2+1;
36     end
37 end
38
39 if c1>c2
40     test_set.gender(i)=1;
41 elseif c2>c1
42     test_set.gender(i)=2;
43 end
44
45 end
46
47 p1 = test_set(test_set.gender==1, :);
48 p2 = test_set(test_set.gender==2, :);
49
50 plot(p1.height, p1.weight, 'c^')
51 plot(p2.height, p2.weight, 'rs')
52
```

Ilustración 7. Análisis de datos.

Antes de la toma de resultados, se le ha dado mayor fiabilidad al algoritmo mediante la comprobación con el modelo ideado para los cálculos analíticos, este se presenta en la .

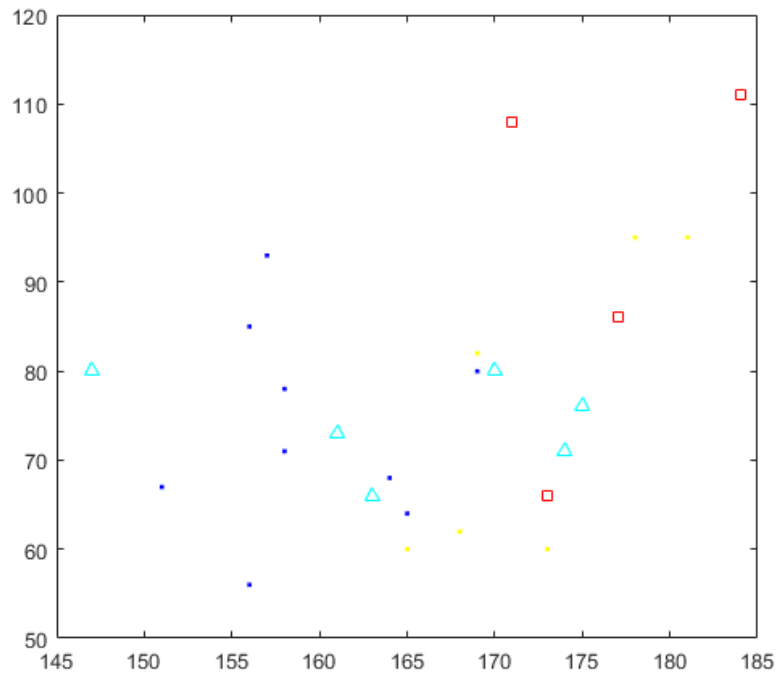


Ilustración 8. Prueba de fiabilidad de código a implementar.

R1.1 Resultados:

Antes de presentar los datos obtenidos para cada valor de k , se presenta en la Ilustración 9 el dataset de entrenamiento.

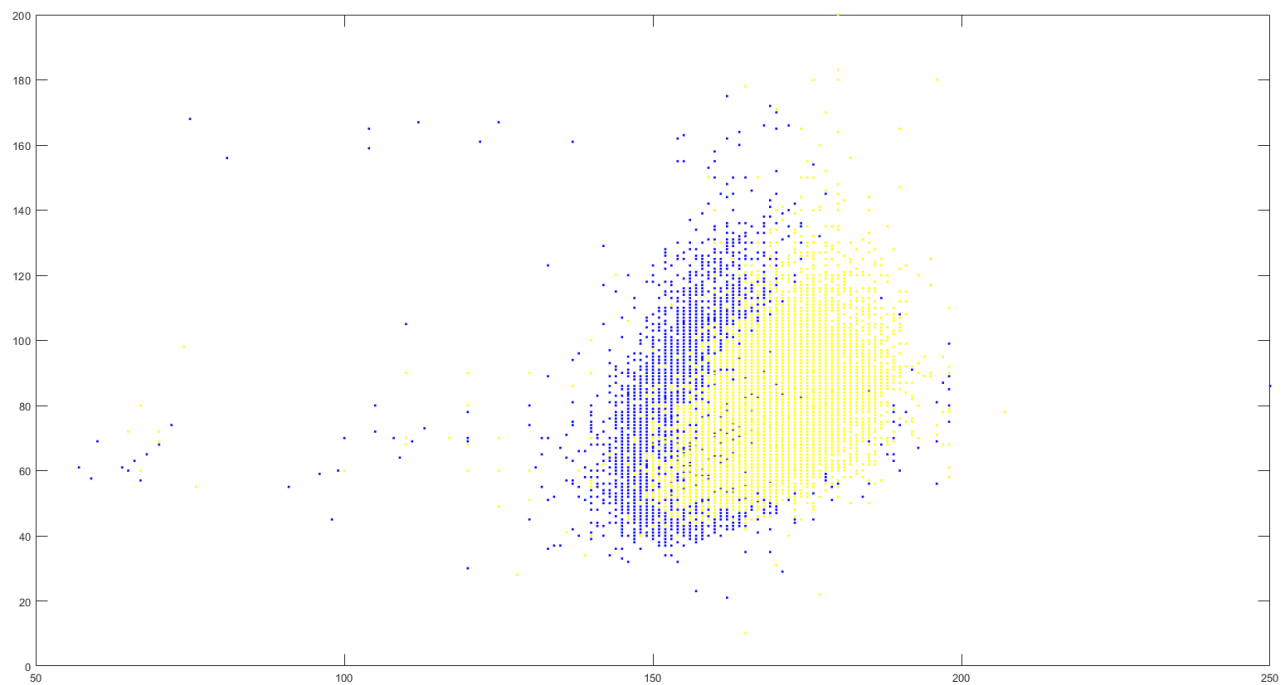


Ilustración 9. Datos de entrenamiento.

A continuación, se presentan todos los gráficos obtenidos. Empezando con $k = 1$ (ver Ilustración 10).

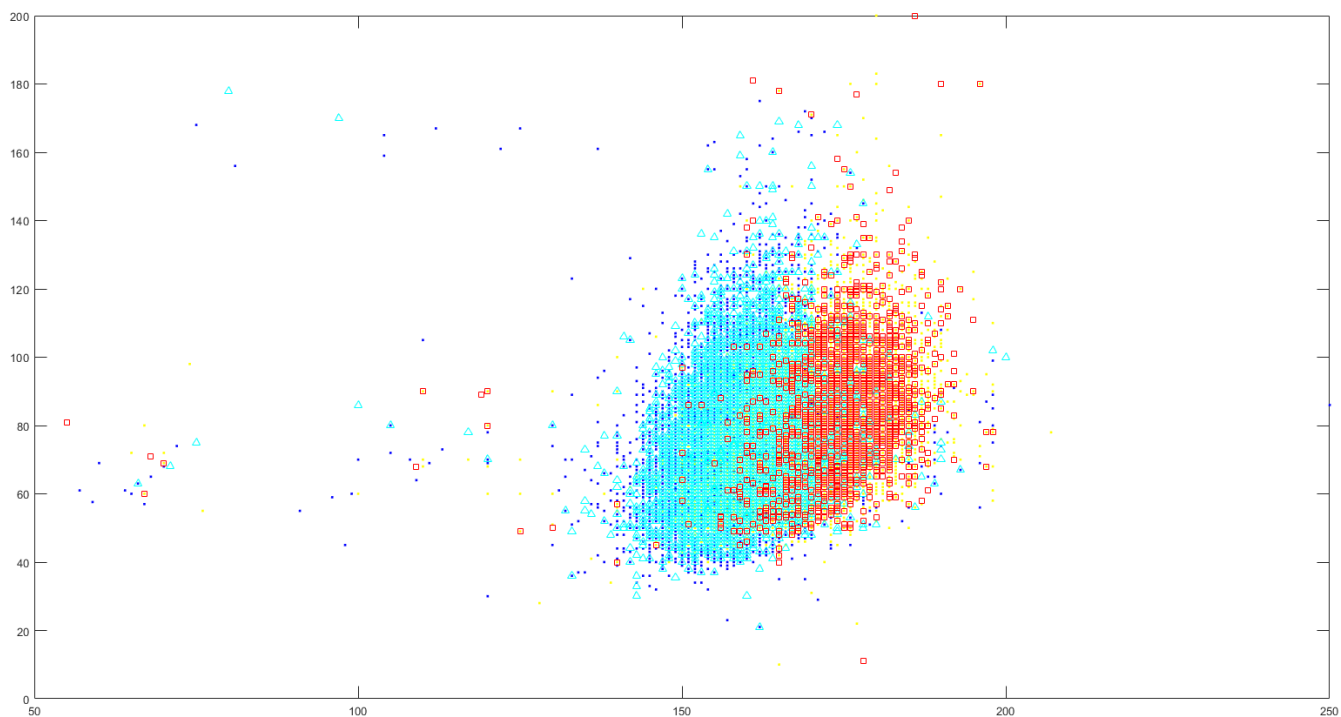


Ilustración 10. resultados para $k=1$.

Para $k = 3$ (ver Ilustración 11).

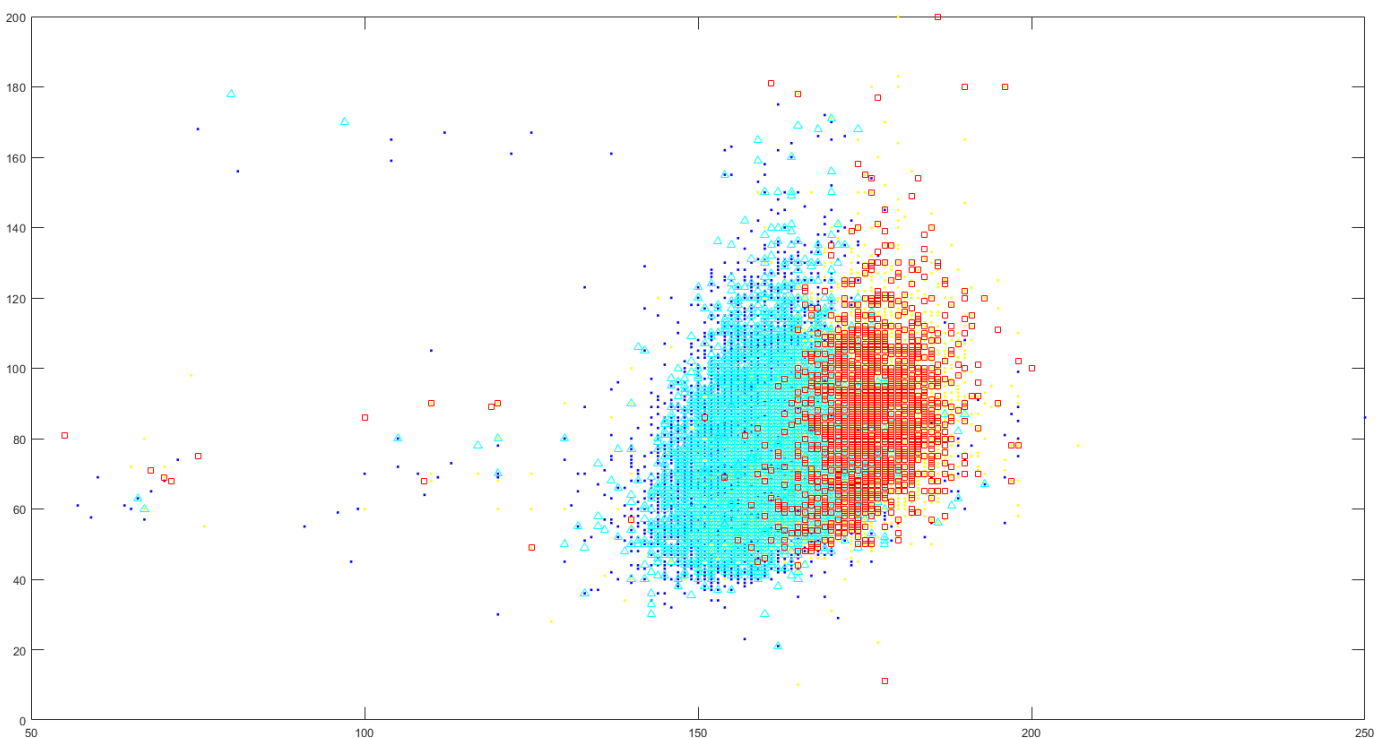


Ilustración 11. Resultados para $k=3$.

Para $k = 5$ (ver Ilustración 12).

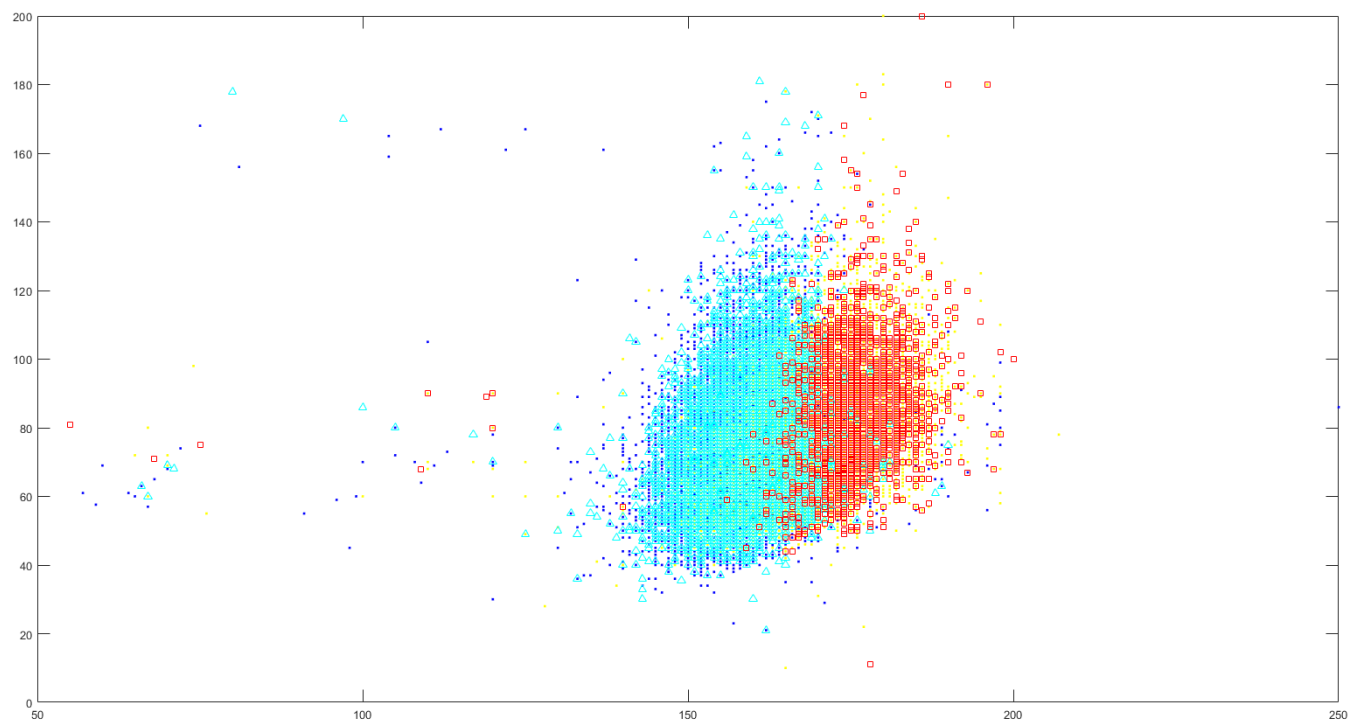


Ilustración 12. Resultados para $k=5$.

Para $k = 9$ (ver Ilustración 13).

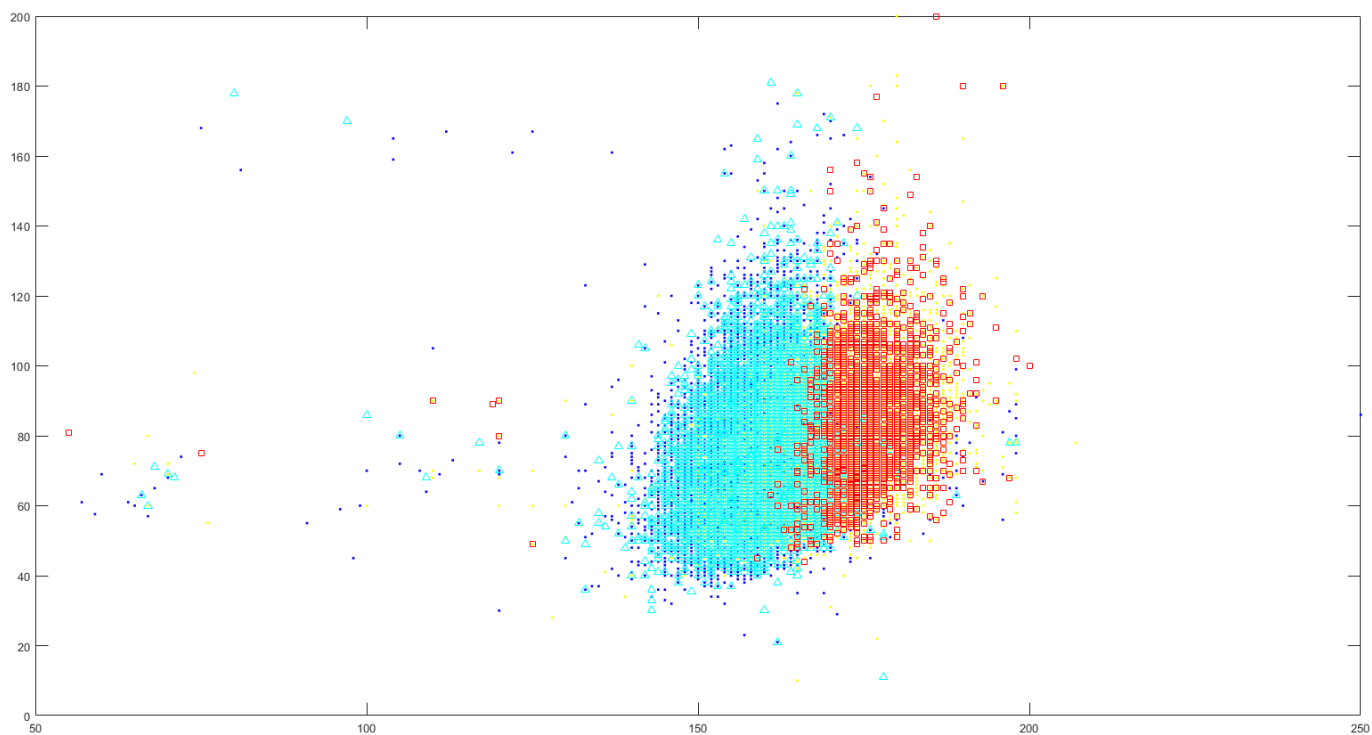


Ilustración 13. Resultados para $k=9$.

Para $k = 15$ (ver Ilustración 14).

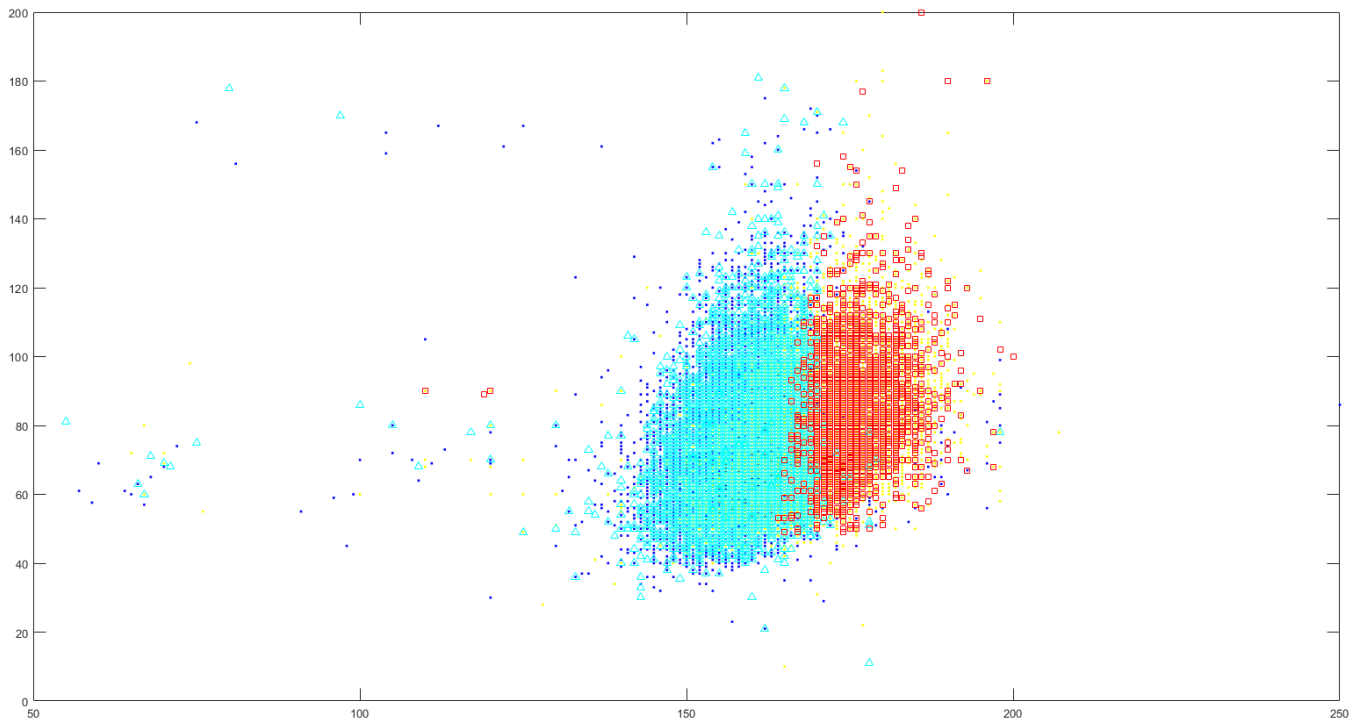


Ilustración 14. Resultados para $k=15$.

A.1.1 Análisis:

Se ha podido observar el comportamiento del algoritmo de los k vecinos más cercanos, proceso de clasificación mediante el cual se hace una comparación entre los datos a ser analizados y su posición en el plano, tomando en consideración los datos ya reconocidos y clasificados que se encuentren en su entorno.

El algoritmo implementado ha sido probado mediante varios valores de k , de donde se interpreta que, para valores de k mayores, tiende a tenerse mayor nitidez en los cúmulos de un valor determinado. Confirmando este algoritmo con la clasificación de los datos calculados analíticamente para tener garantía de su correcto funcionamiento. Se ha seleccionado para todo caso una distancia euclídeana.