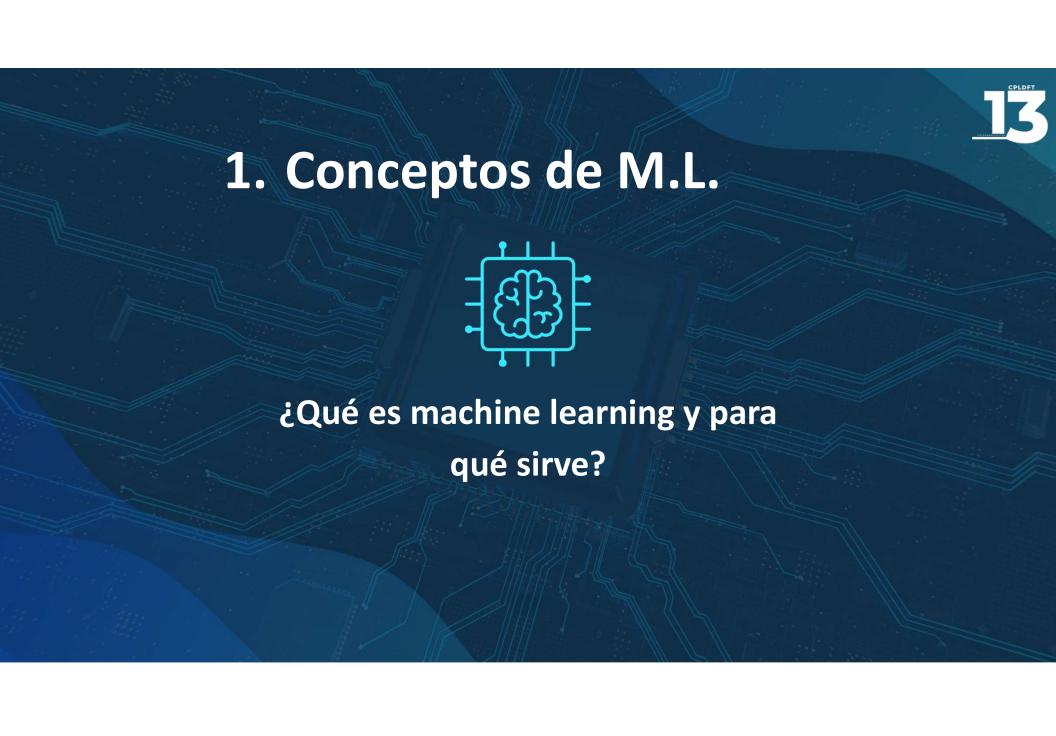


Disclaimer

- El taller **no** busca convertir a los participantes en científicos de datos.
- Se busca, a través de ejemplos sencillos, que los participantes puedan comprender los conceptos fundamentales de IA y ML y ver su aplicabilidad en monitoreo transaccional.
- Se busca que, un OC pueda entender cómo se aplican los modelos, entienda las salidas y pueda tener una conversación con los equipos de ML de su organización.





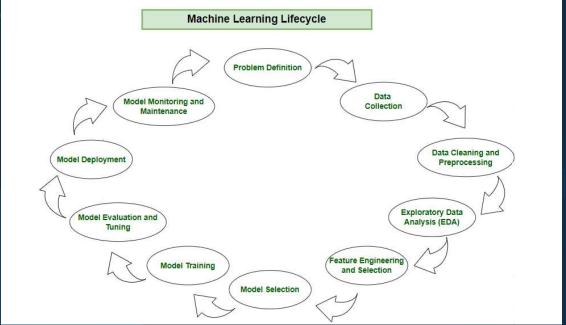


Definición técnica

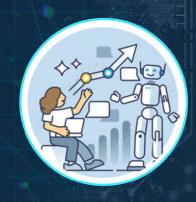
Machine Learning es una rama de la inteligencia artificial que desarrolla algoritmos capaces de aprender de los datos para reconocer patrones, identificar anomalías y realizar predicciones de manera automática.



El ciclo de implementación



Modelos de aprendizaje automático



Supervisado

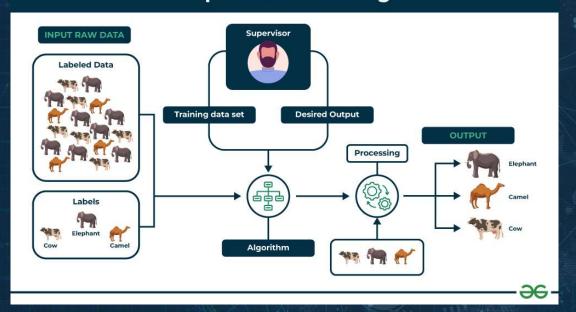


No Supervisado



Aprendizaje supervisado

Supervised Learning

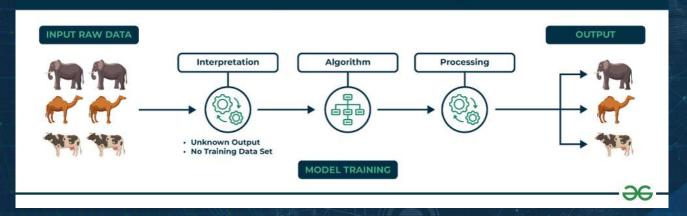






Aprendizaje no supervisado

Unsupervised Learning





¿Qué algoritmos usamos para M.L.?

Regresión Lineal

Decision Tree

Random Forest

K-Means

DBSCAN

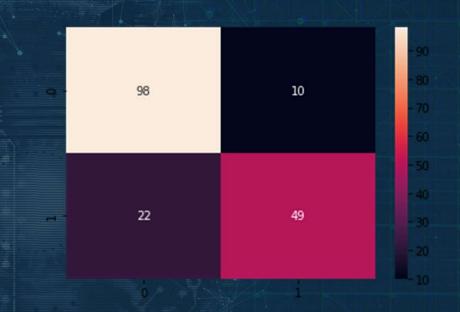
Gaussian Mixture Model



Ejemplo: Supervivencia del Titanic



https://github.com/alicevillar/titanickaggle/blob/main/Titanic DecisionTree.ipynb



b) Precision

Precision = true positive / true positive + false positive precision_dt = metrics.precision_score(Y_real,Y_pred_dt) precision_dt

0.8305084745762712

La pregunta es fácil

¿Por qué no creamos un modelo de M.L. con la capacidad de medir la probabilidad de que un perfil o comportamiento esten ligados a lavado de activos?

La respuesta no es igual de fácil





Limitaciones

- Base de Datos no balanceada
- El impacto de los falsos negativos
- Falta de datos para un "labeling" real



Los datos claves no siempre están



La Buena Noticia

Que no nos podamos saltar al final no significa que no nos pueda ayudar durante el proceso

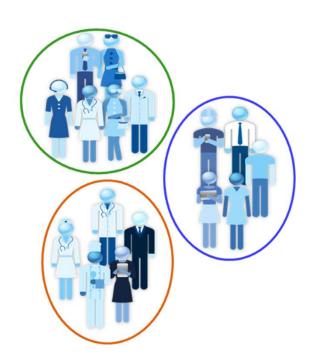




2. Uso de cluster como herramienta de monitoreo

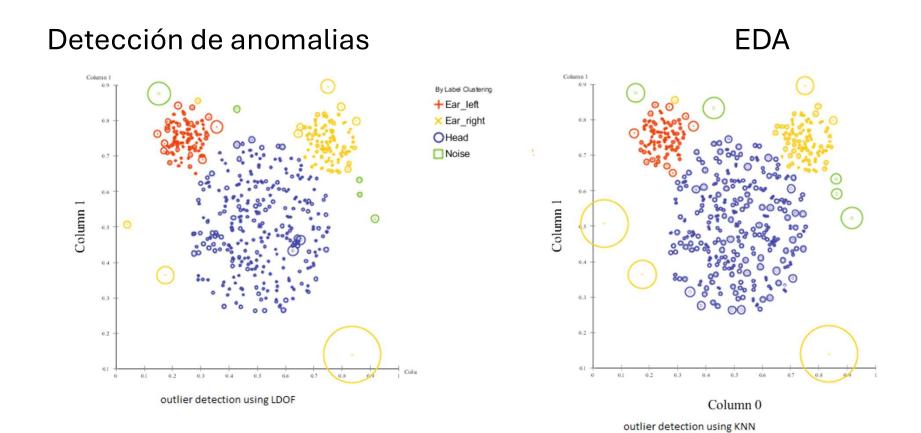
Qué es un cluster?

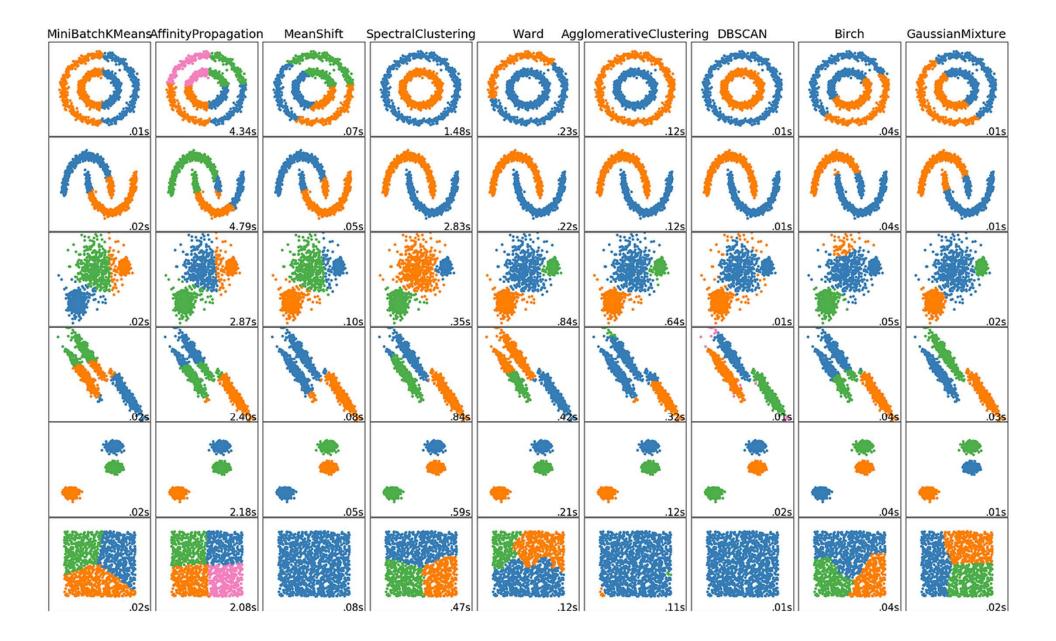




• Técnica que busca encontrar grupos (clusters) de puntos de datos de tal forma que los puntos de un mismo grupo sean similares (o parecidos) entre si y al mismo tiempo, diferentes (o no relacionados) con puntos de datos en otros grupos.

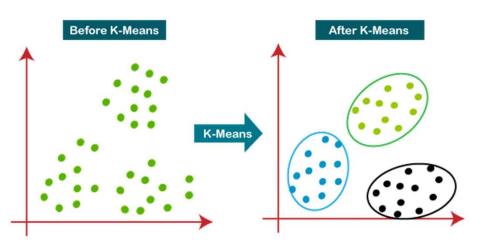
Otros usos de cluster



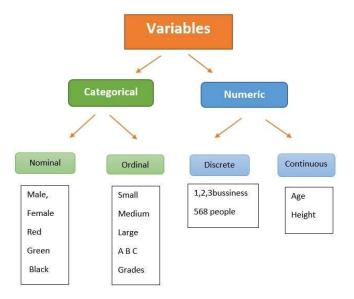


Kmeans cluster - Requisitos

Grupos hiperesféricos



Variables numéricas



Ejercicio

- Tenemos la siguiente data de KYC y deseamos disenar un modelo de cluster para segmentar clientes desde la perspectiva de riesgo AML. Cuáles variables podríamos utilizar?
 - Nombre
 - Direction
 - Teléfono
 - Tipo (Natural, Jurídico)
 - Nacionalidad
 - Actividad Ecónomica
 - Ocupación
 - Nivel de estudios

- Productos / servicios
- Canales que utiliza
- Vol Txs money in
- Q Txs money in
- Vol Txs money out
- Q Txs money out

Ejercicio

• Más variables no siempre es mejor. Cuáles podriamos "agrupar"?

- Nombre
- Direction
- Teléfono
- Tipo (Natural, Jurídico)
- Nacionalidad
- Actividad Ecónomica
- Ocupación
- Nivel de estudios

- Productos / servicios
- Canales que utiliza
- Vol Txs money in
- Q Txs money in
- Vol Txs money out
- Q Txs money out

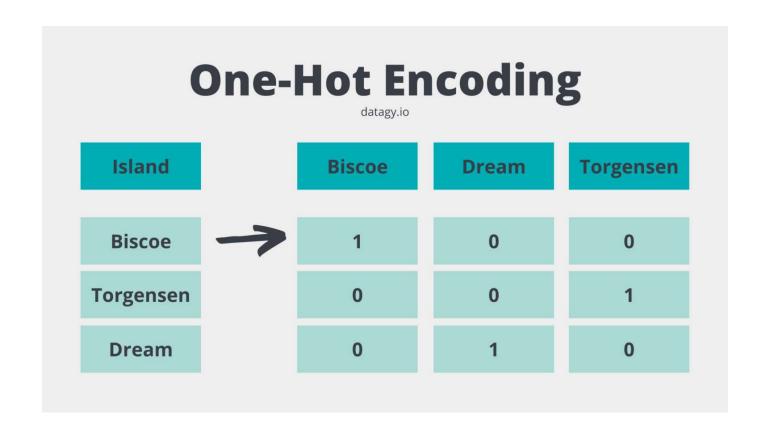
Z Score Formula



$$z = \frac{\lambda - \mu}{\sigma}$$



Kmeans: Transformacion de datos









• Se tiene la siguiente base de datos, cómo se podría transformar?

IdCliente	NR LAFT	Nacional	PEP	Vol Txs	Q Txs	Vol Txs	Q Txs money	
				money in	money in	money out	out	
1	Alto	Nacional	Si	12000	20	11500	14	
2	Medio	Nacional	No	132000	67	113450	56	
3	Medio	Extranjero	No	6400	45	2400	32	
4	Bajo	Nacional	No	500	28	500	16	
5	Alto	Nacional	No	276000	3	267000	1	

Kmedias: Ventajas y desventajas

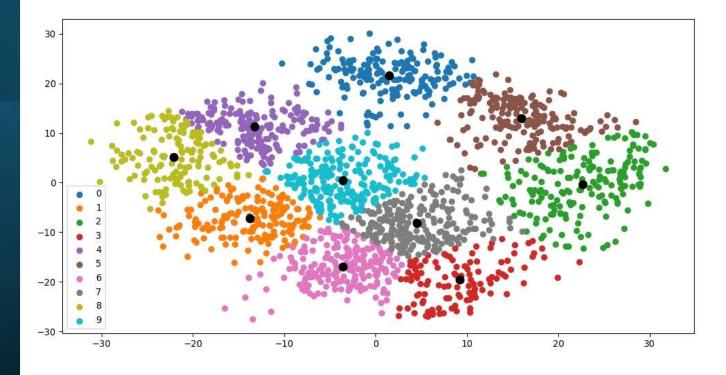
Ventajas

- Simple y fácil de entender
- Rápido y escalable (útil en datasets grandes)
- Puede manejar diferentes tipos de datos*
- Puede ser utilizado en una amplia gama de aplicaciones

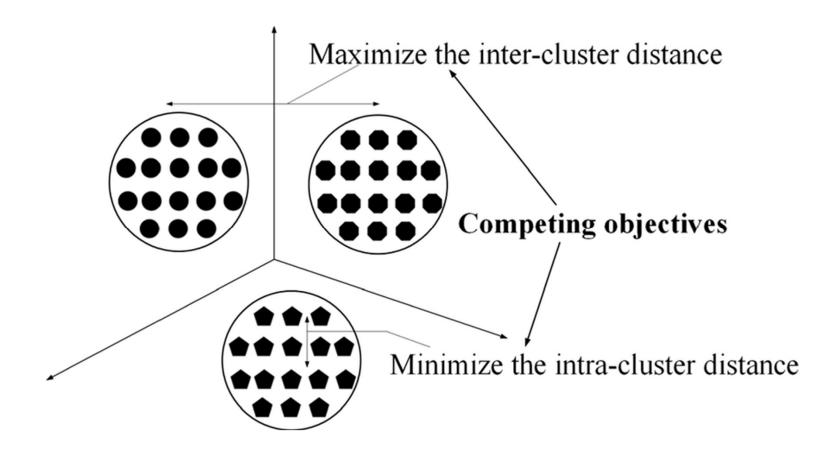
Limitantes

- Requiere que se especifique el numero de segmentos con antelación.
- Sensible a la inicialización de centroides
- Problemas con data con estructura compleja
- Puede verse afectado por outliers

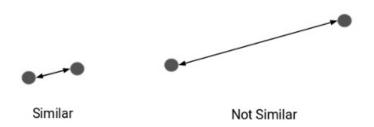
Kmedias y centroides



Cómo se mide la similaridad?



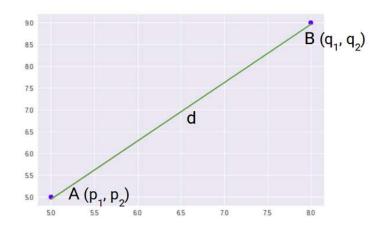
Similaridad (homogeneidad)



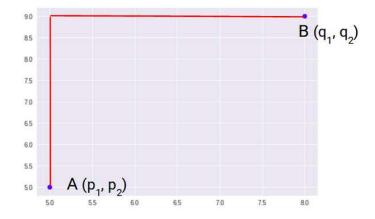
Minkowski

$$D\left(X,Y
ight) = \left(\sum_{i=1}^{n}\left|x_{i}-y_{i}
ight|^{p}
ight)^{1/p}$$

Euclidiana



Manhattan





Euclidean distance between pairs of samples														
	1	0	0.1129	0.2283	0.29	0.4016	0.5018	0.4774	0.09234	0.1143	0.1382			
	2	0.1129	0	0.1833	0.1932	0.3372	0.4053	0.376	0.149	0.05639	0.1916			0.5
	3	0.2283	0.1833	0	0.1604	0.1764	0.3375	0.3541	0.2982	0.1737	0.3467			0.4
No.	4	0.29	0.1932	0.1604	0	0.2093	0.223	0.2146	0.3294	0.1925	0.3707			
Second Sample No.	5	0.4016	0.3372	0.1764	0.2093	0	0.2668	0.3272	0.4666	0.3325	0.5156		-	0.3
S puoc	6	0.5018	0.4053	0.3375	0.223	0.2668	0	0.1022	0.5326	0.3994	0.5657			
Se	7	0.4774	0.376	0.3541	0.2146	0.3272	0.1022	0	0.4942	0.3724	0.5204		,_	0.2
	8	0.09234	0.149	0.2982	0.3294	0.4666	0.5326	0.4942	0	0.1516	0.07304			
	9	0.1143	0.05639	0.1737	0.1925	0.3325	0.3994	0.3724	0.1516	0	0.1981		-	0.1
1	0	0.1382	0.1916	0.3467	0.3707	0.5156	0.5657	0.5204	0.07304	0.1981	0			0
1 2 3 4 5 6 7 8 9 10 First Sample No.												U		

3. Análisis preliminar

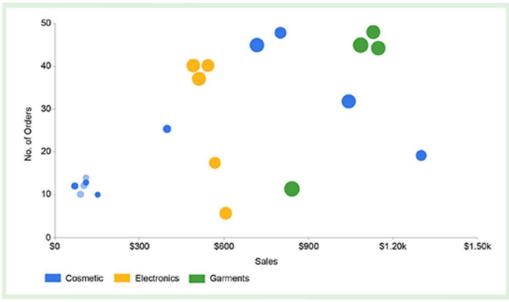
Matriz de correlaciones

Exploratory Data Analysis

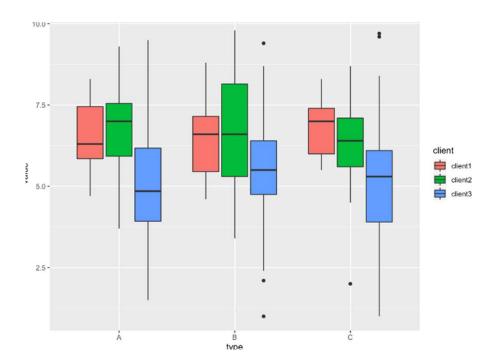
goai		0.73	0.7	0.40	0.21	-0.21	-0.50	-0.50	-0.43	-0.10	0.21
am	0.79	1	0.71	0.6	0.17	-0.23	-0.69	-0.59	-0.52	-0.24	0.06
drat	0.7	0.71	1	0.68	0.44	0.09	-0.71	-0.71	-0.7	-0.45	-0.09
mpg	0.48	0.6	0.68	1	0.66	0.42	-0.87	-0.85	-0.85	-0.78	-0.55
VS	0.21	0.17	0.44	0.66	1	0.74	-0.55	-0.71	-0.81	-0.72	-0.57
qsec	-0.21	-0.23	0.09	0.42	0.74	1	-0.17	-0.43	-0.59	-0.71	-0.66
wt	-0.58	-0.69	-0.71	-0.87	-0.55	-0.17	1	0.89	0.78	0.66	0.43
disp	-0.56	-0.59	-0.71	-0.85	-0.71	-0.43	0.89	1	0.9	0.79	0.39
cyl	-0.49	-0.52	-0.7	-0.85	-0.81	-0.59	0.78	0.9	1	0.83	0.53
hp	-0.13	-0.24	-0.45	-0.78	-0.72	-0.71	0.66	0.79	0.83	1	0.75
carb	0.27	0.06	-0.09	-0.55	-0.57	-0.66	0.43	0.39	0.53	0.75	1

Técnicas de reducción de dimensiones

Qué pasa cuando se tienen más de 3 variables? Cómo visualizarlas en conjunto?







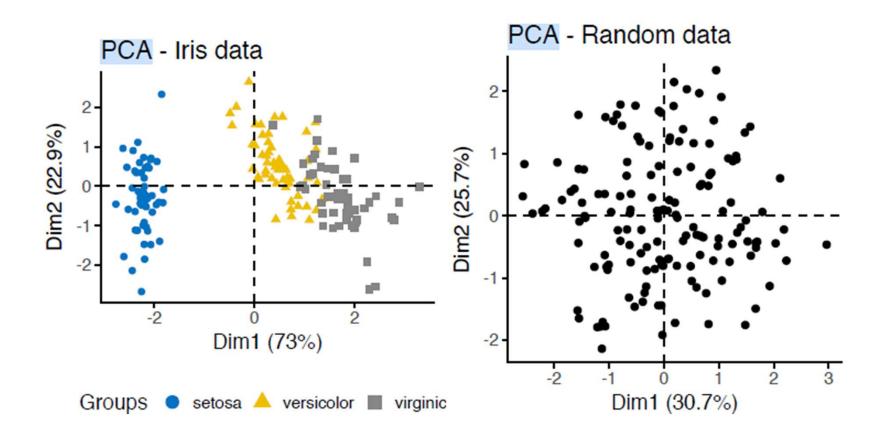
How to Create a Scatter Plot in Excel With 3 Variables?

PCA

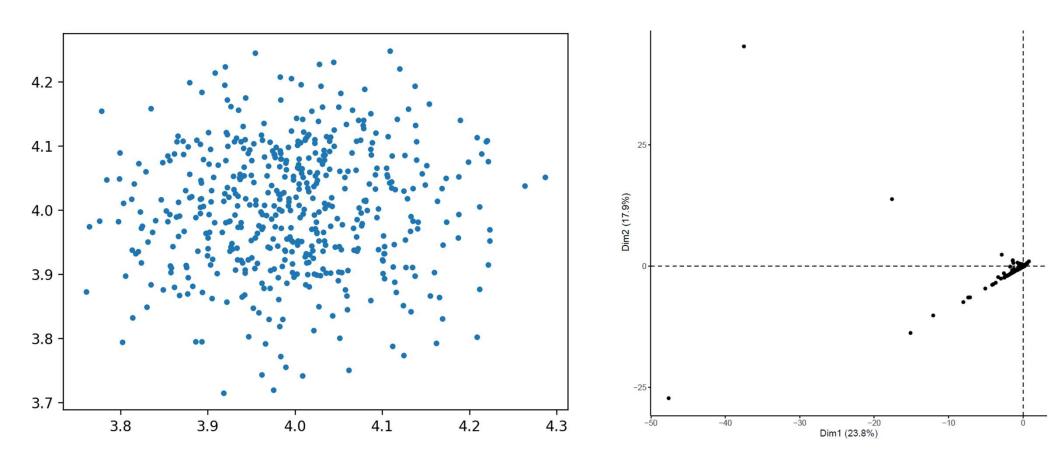
- Identifica "direcciones" que maximizan varianza
- Sobre ejes ortogonales entre si (correlacion 0)
- Las "direcciones" o nuevas variables, son una combinación lineal de las variables originales
- Retorna la "mejor reconstrucción", es decir, minimiza la Perdida de informacion al moverse de N (dimensiones PCA) a M (dimensiones originales)
- Los valores propios son monotonicamente decrecientes (se reducen en cada nueva dimension) y representan el nivel de varianza que acumulan.
- Bien estudiado -> algoritmos rapidos



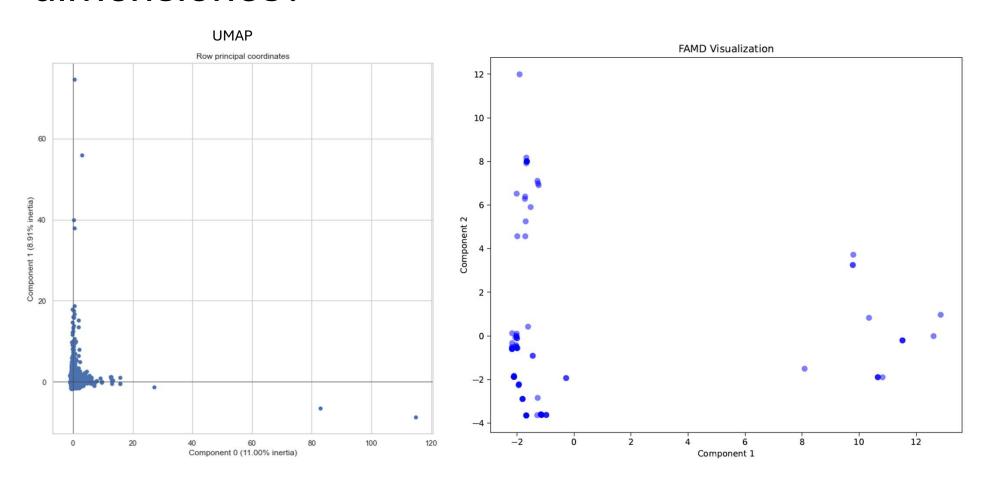
PCA: Inspección visual de los datos para determinar "clustereabilidad"



PCA: Ejercicio – Se podría aplicar cluster?

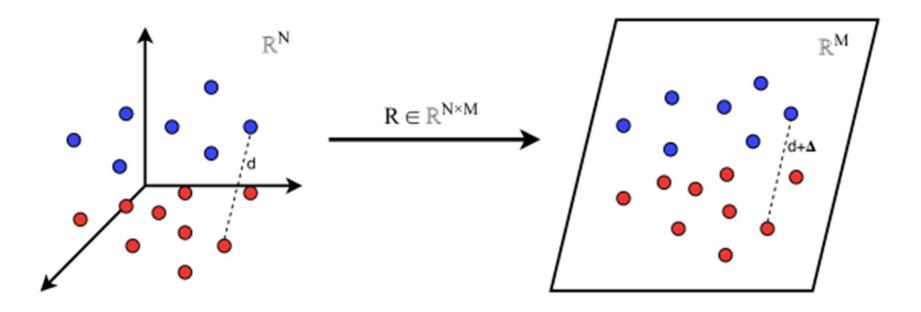


PCA es la unica técnica de reduccion de dimensiones?



Porqué no montar el modelos de cluster directamente en los componentes principales?

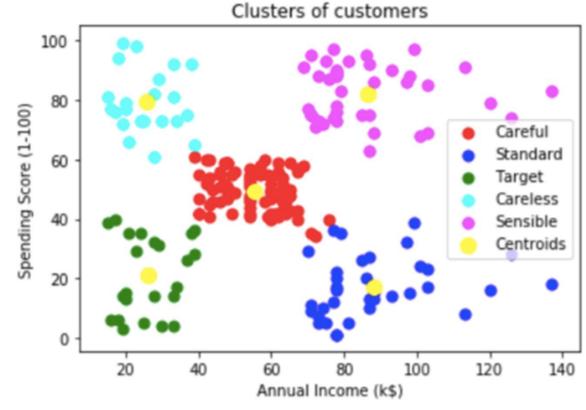
- Interpretabilidad
- Expectativa del regulador
- "Perdida de informacion"

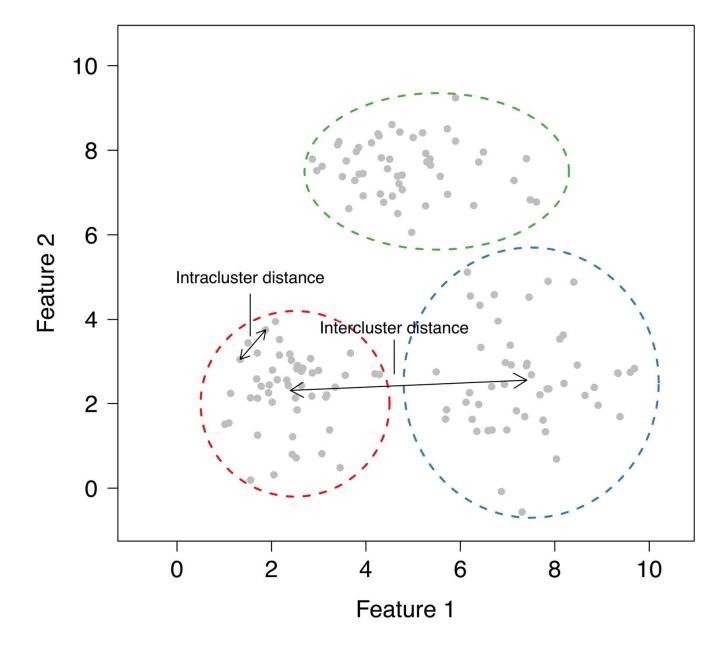


4. Selección del número óptimo de segmentos

Porque seleccionar un número de segmentos?

- El parámetro k es un input del modelo
- De su elección dependen aspectos críticos, p.e., definición de las senales de alerta.





4.a El criterio del codo (Elbow criterion)



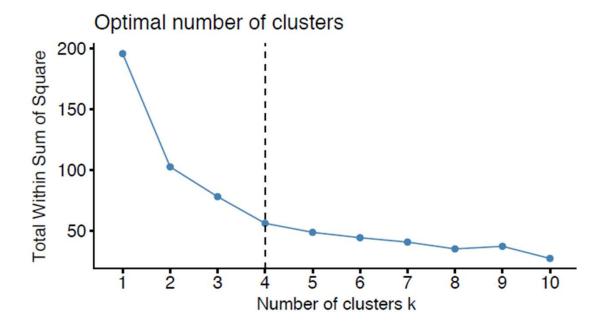
Criterio del codo

Importancia del centroide

$$W(C_k) = \sum_{x_i \in C_k} (x_i - \mu_k)^2$$

 x_i design a data point belonging to the cluster C_k μ_k is the mean value of the points assigned to the cluster C_k

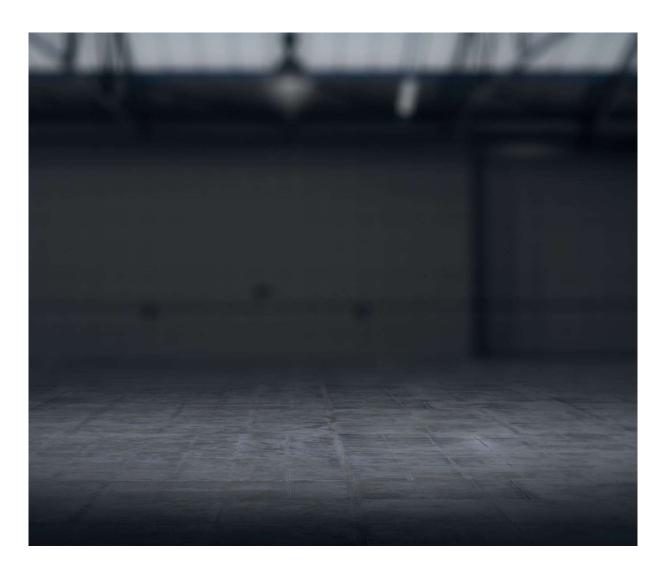
ot.withinss =
$$\sum_{k=1}^{k} W(C_k) = \sum_{k=1}^{k} \sum_{x_i \in C_k} (x_i - \mu_k)^2$$

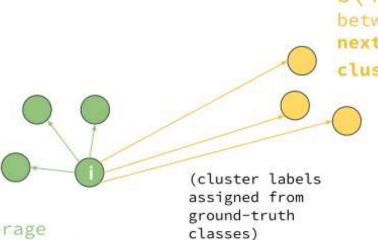












a(i): average
distance between i
and all of the
other points in its
own cluster

b(i): distance between i and its next nearest cluster centroid

For a single point, i

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i),b(i)\}}$$

$$-1 \leq s(i) \leq 1$$
 Sprawling, overlapped clusters Tight, well-separated clusters

Silueta

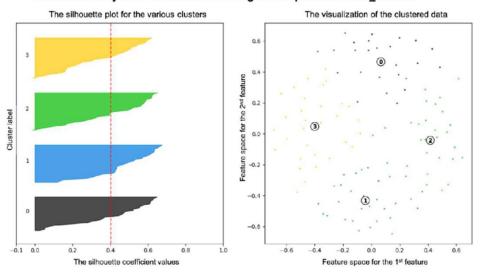
Para cada punto se calcula:

- La distancia promedio a todos los demas puntos pertenecientes al mismo cluster (cohesion)
- La distancia promedio a todos los puntos pertenecientes al segmento mas cercano (separacion)

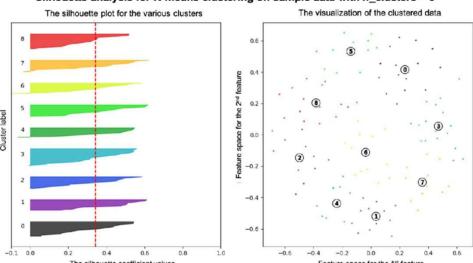
La silueta se calcula como:

 (separacion – cohesion) / max (separacion y cohesion)

Silhouette analysis for K-means clustering on sample data with $n_{clusters} = 4$

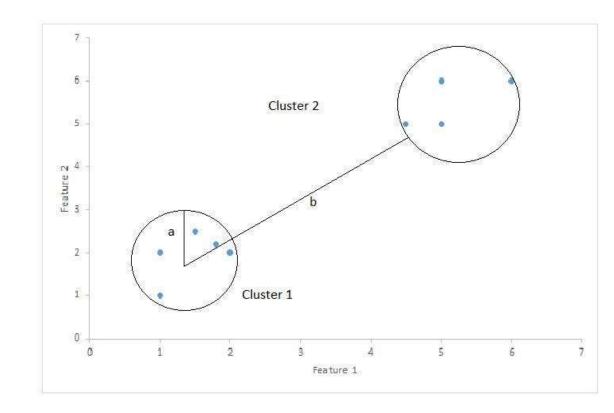


Silhouette analysis for K-means clustering on sample data with $n_{\rm clusters}$ = 9



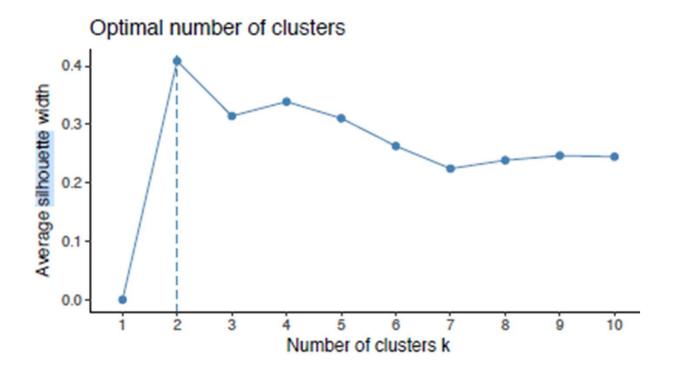
Silueta

- Valores de silueta:
- 1: Los segmentos estan separados entre si y estan claramente diferenciados
- 0: Los segmentos son indiferentes, o que la distancia es insignificante
- -1:La asignacion de segmentos ha sido errónea

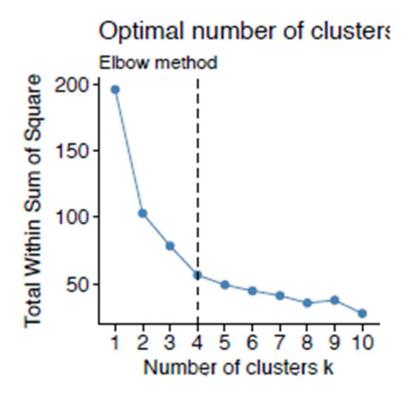


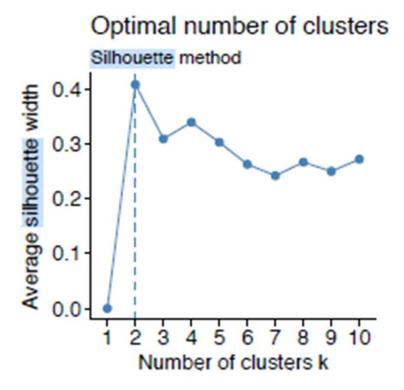
Cual valor de K seleccionar en function de la Silueta?

Silueta

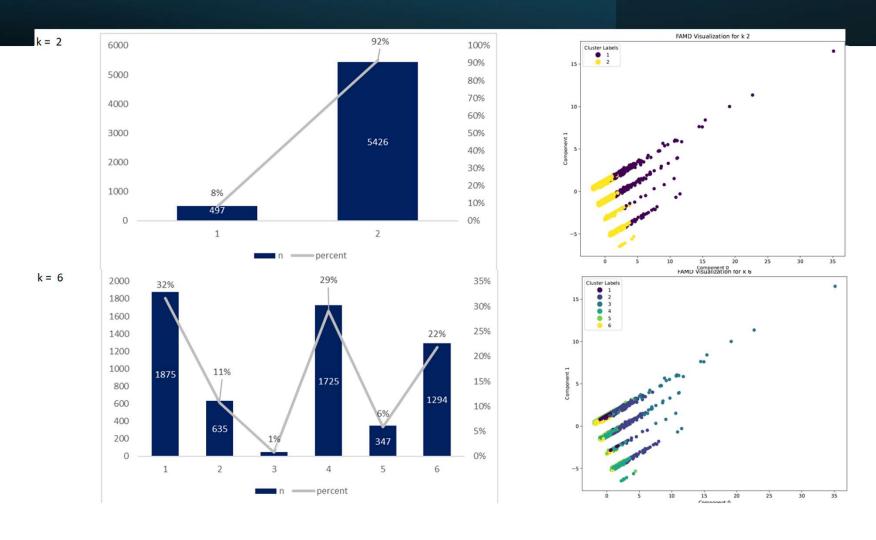


Putting all together

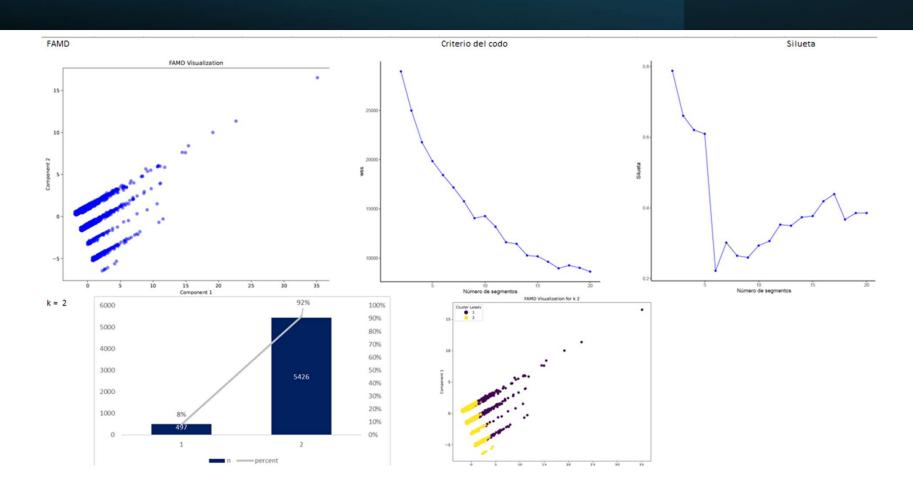




Otros aspectos a evaluar: Concentración de segmentos



Ejercicio 2: Cuál sería su recomendación?

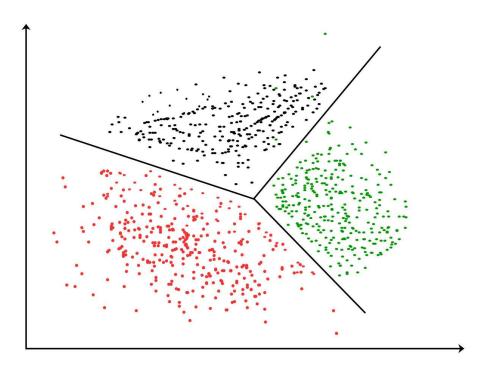




5. Monitoreo

Una vez definido el modelo, como hacemos monitoreo?

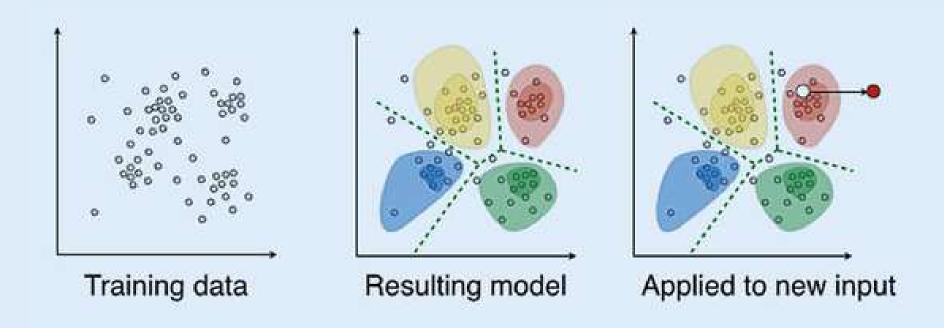
Resegmentar



Modelo de pronostico

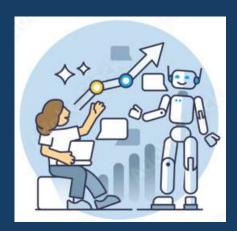


Modelo de pronostico



Bajo que grupo de técnicas se podría diseñar un modelo de pronostico de segmentos

Aprendizaje Supervisado



Data "marcada"

Task Drive

Predecir el siguiente valor

Aprendizaje No Supervisado

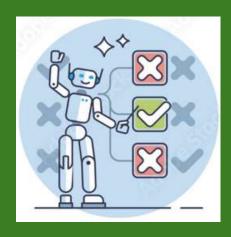


Data "no marcada"

Data Driven

Identificar estructura en los datos

Aprendizaje por refuerzo



Aprender de errores



Interquartile Range & Outliers

Ejercicio: Diseñemos umbrales de alertamiento basados en los resultados de la segmentación

Α	В	С	D	E	F	G	Н	l l	J
dCliente	NRActividad	NR	Nacional	PEP	QTxsCASHIN	QTxsCASHOUT	suma_montoCASHIN	suma_montoCASHOUT	Cluster
1	Medio	Bajo	1	0	5	3	61,38175	59,5	6
2	Alto	Bajo	1	0	42	26	11079,2485	7875,72075	2
3	Alto	Alto	1	0	3	2	3318,05375	3400	4
4	Medio	Bajo	1	0	2	4	4976,6665	2741,0055	6
5	Alto	Bajo	1	0	27	22	293,23925	292,685	1
6	Alto	Bajo	1	0	1	0	0,05	0	1
7	Alto	Bajo	1	0	3	4	95,6075	95,8605	1
8	Alto	Alto	1	0	4	4	2036,7225	3195,677	4
9	Alto	Alto	1	0	5	5	455,12425	458,6805	4
10	Alto	Alto	1	0	3	2	60,55775	115,192	4
11	Alto	Bajo	1	0	2	3	156,14475	155,5	1
12	Bajo	Medio	1	0	6	2	215,96025	210	6
13	Alto	Bajo	1	0	81	0	255,22325	0	1
14	Alto	Bajo	1	0	2	1	64,58475	64	1
15	Alto	Medio	1	0	1	1	2,6705	2,5	1
16	Alto	Alto	0	0	1	0	1,11275	0	4
17	Alto	Bajo	1	0	15	20	2370,29725	2358,24425	1
18	Alto	Alto	1	0	3	9	414,1965	407,4345	4





Brian S. Everitt • Sabine Landau Morven Leese • Daniel Stahl

Springer Series in Statistics

Trevor Hastie Robert Tibshirani Jerome Friedman

The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Second Edition

Multivariate Analysis I

Alboukadel Kassambara

Practical Guide To Cluster Analysis in R

Unsupervised Machine Learning

sthda.com Edition 1

Bibliografia