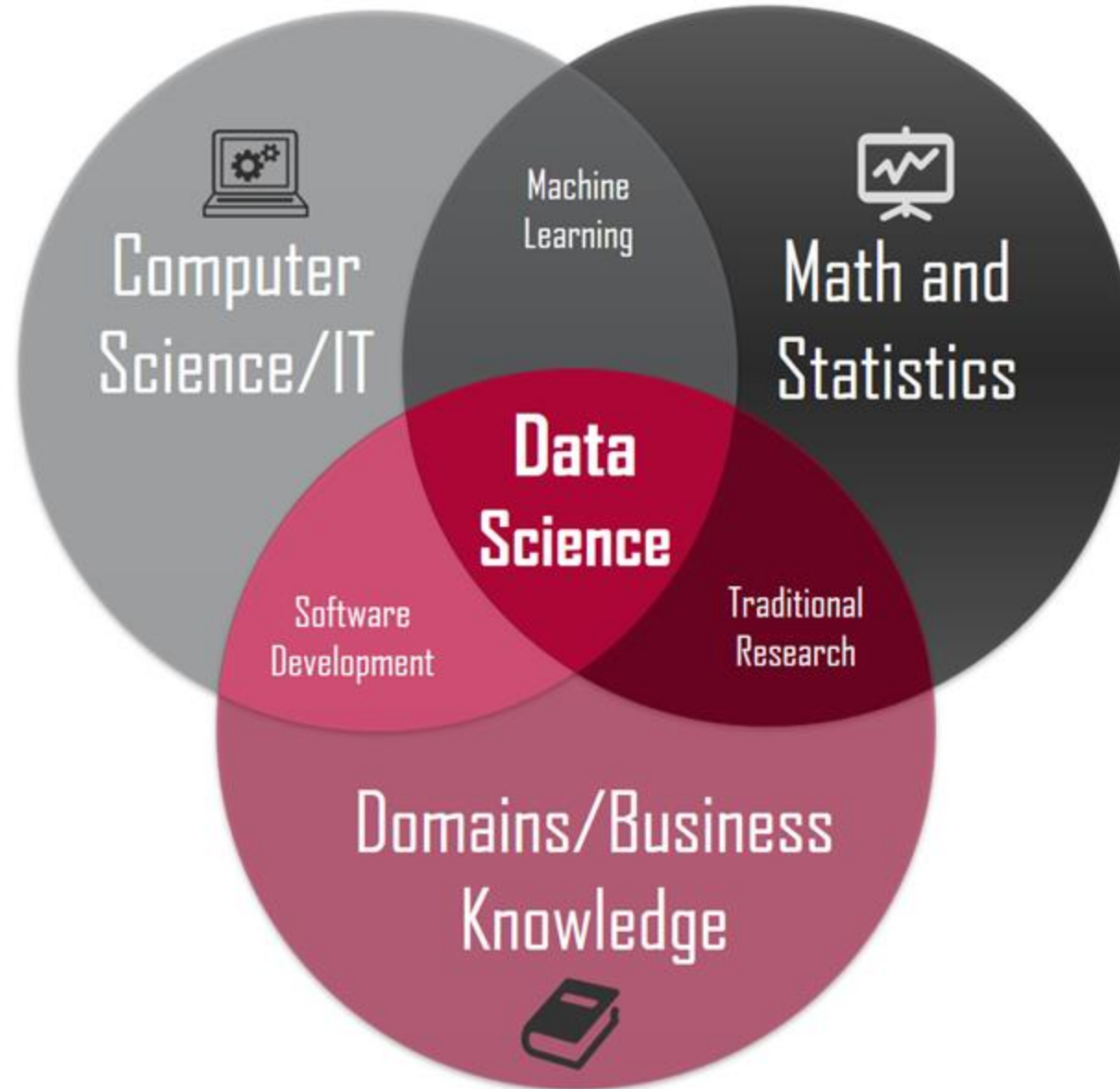


Introducción a la Ciencia de Datos

Sesión 1

¿Qué herramientas utiliza un Científico de Datos?



¿Qué lenguaje de programación utilizar?



VS.



¿Dónde comenzar a aprender rápido y gratis?

The Kaggle logo, featuring the word "kaggle" in a light blue, lowercase, sans-serif font.

<https://www.kaggle.com/learn/overview>



<https://courses.analyticsvidhya.com/>

The Towards Data Science logo, featuring the text "Towards Data Science" in white, serif font, with a horizontal line underneath, all set against a dark blue rectangular background.

<https://towardsdatascience.com/>

METODOLOGÍA CRISP - DM

Cross-Industry Standard Process for Data Mining (finales 90's)

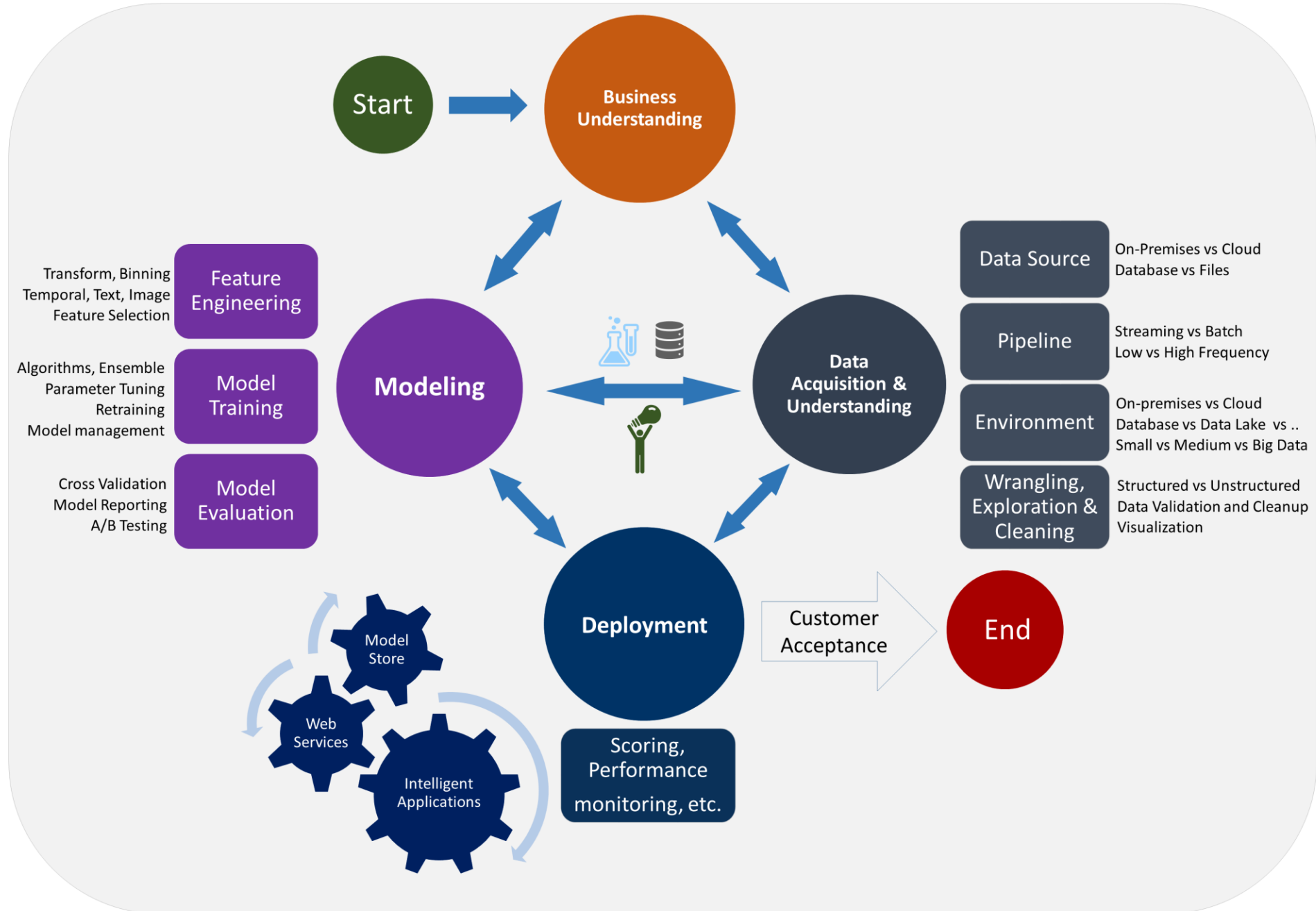


Data Science Lifecycle

**METODOLOGÍA
TDSP**

**Team
Data
Science
Process**

(Microsoft, 2016)



**Identifica una
necesidad del
mercado**

1

¿Algún problema
que puedas
resolver aplicando
análisis de datos?

**Captura los
datos
necesarios y
limpia la
información**

2

¿Tienes todos los
datos que
necesitas? ¿Los
datos están
legibles?

**Selecciona y
evalúa el
modelo**

3

¿El modelo es el
mejor para los
datos que tienes?
¿El modelo es
preciso?

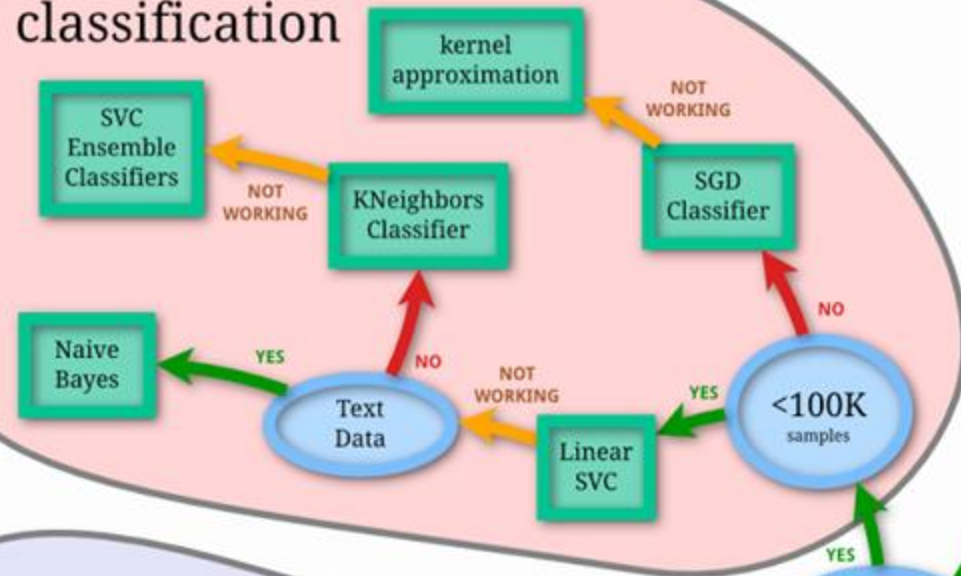
**Implementa el
modelo para
uso de tus
clientes**

4

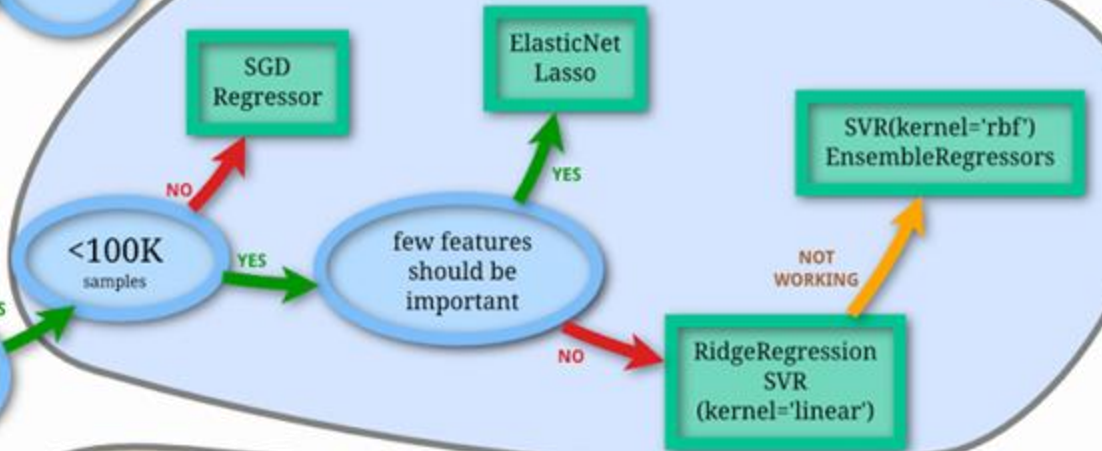
¿Tu modelo puede
salir a producción?
¿Tienes la
infraestructura?

scikit-learn algorithm cheat-sheet

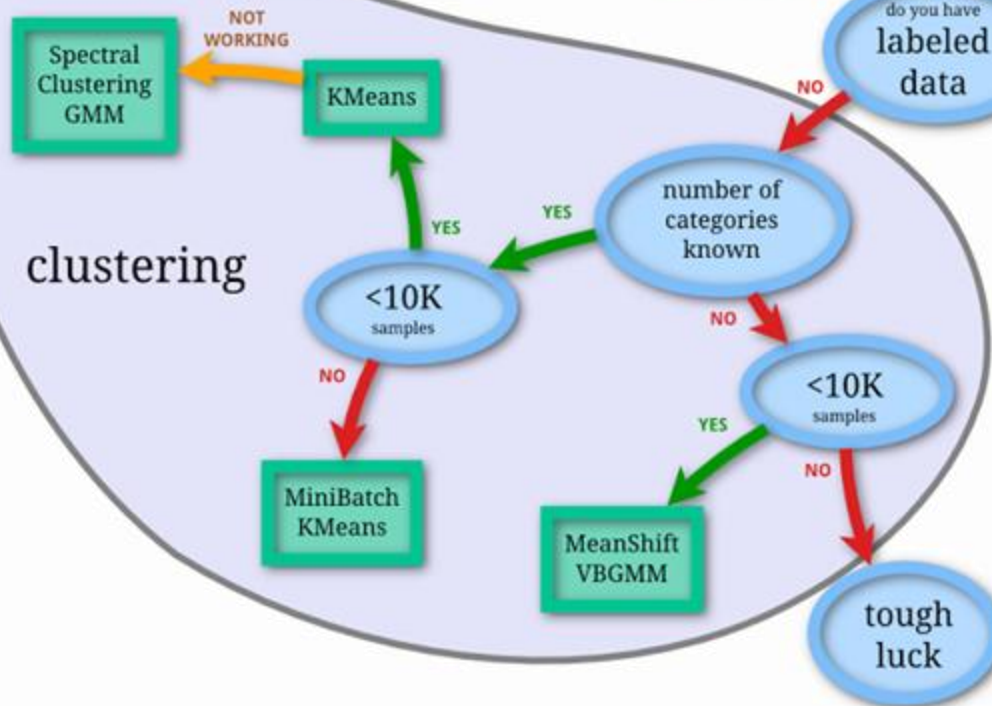
classification



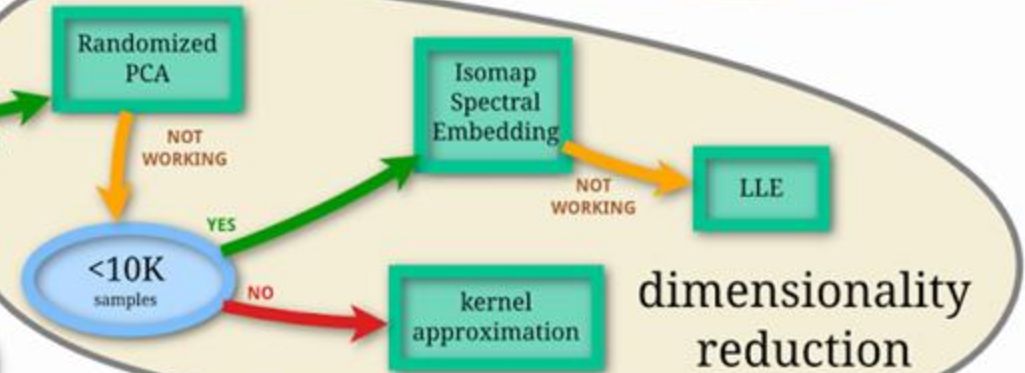
regression



clustering



dimensionality reduction



Aprendizaje No Supervisado

Clusterización con Algoritmo Kmeans

AGRUPAR ELEMENTOS



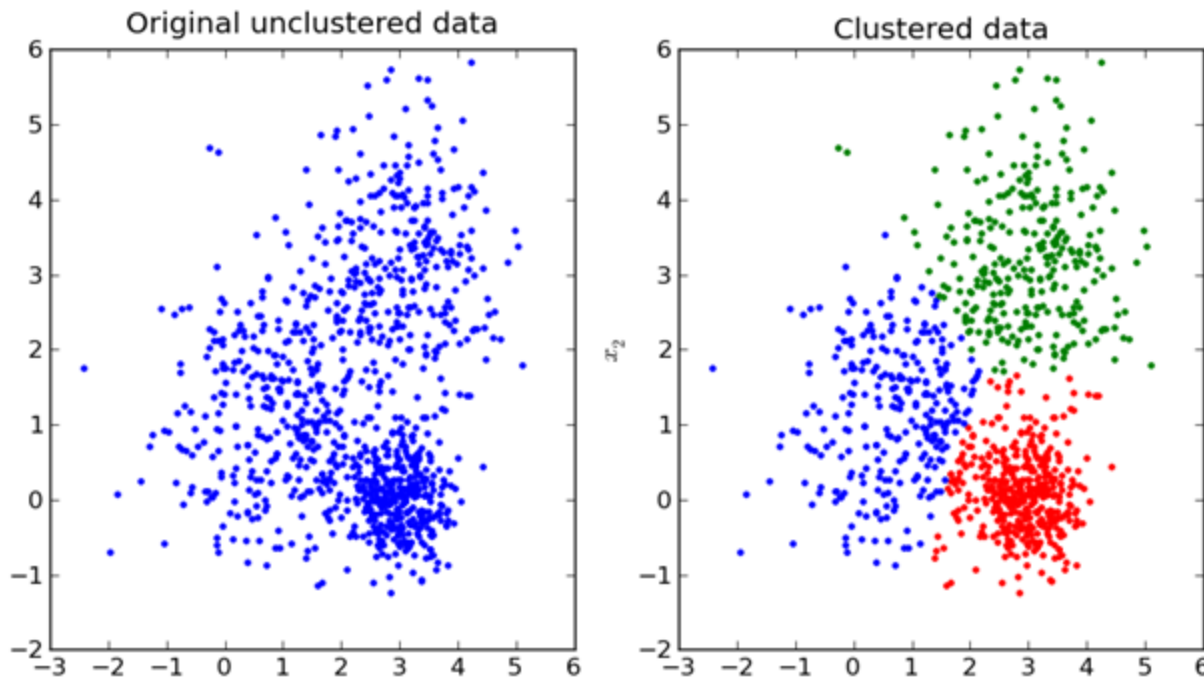
- ✓ Recopilar un conjunto
- ✓ Escoger un criterio aleatorio
- ✓ Ejecutar criterio sobre el conjunto

ORDENAR Y CLASIFICAR



- ✓ Identificar cada elemento
- ✓ Descubrir una característica en común
- ✓ Dar escala a esa característica
- ✓ Asignar cada elemento a un grupo según escala

CLUSTERING



1. Recopilar datos
2. Identificar centroides
3. Asignar cada dato a un *centroide*¹
4. Formar clusters

¹<https://youtu.be/V28U9cGeSdA>

FUNDAMENTO MATEMÁTICO

$d(\mathbf{p}, \mathbf{q}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2}$. ✓ Definir una distancia **euclidiana**

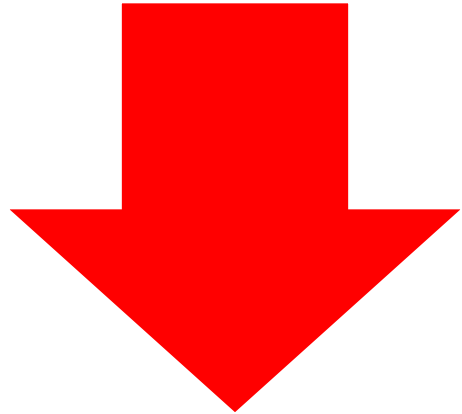
$$\arg \min_{c_i \in C} \text{dist}(c_i, x)^2$$

✓ Asignar cada dato al centroide (c_i) separado por la mínima distancia

$$c_i = \frac{1}{|S_i|} \sum_{x_i \in S_i} x_i$$

✓ Hallar nuevo centroide y repetir

EQUILIBRIO BUSCADO



Minimizar la
distancia entre
cada elemento
dentro del
cluster)

✓ Los elementos de un mismo grupo son muy similares



Maximizar la
distancia entre
cada elemento
fuera del
cluster



✓ Un grupo se distingue de otro fácilmente

REQUISITOS DEL MODELO

- ✓ Los datos deben ser numéricos y continuos (no flags)

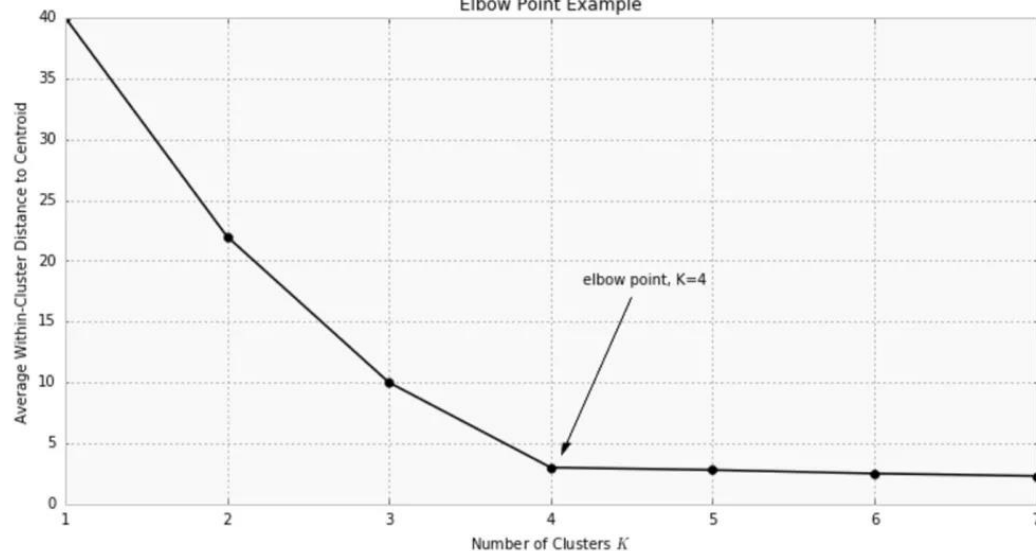
	usuario	op	co	ex	ag
0	3gerardpique	34.297953	28.148819	41.948819	29.370315
1	aguerosergiokun	44.986842	20.525865	37.938947	24.279098
2	albertochicote	41.733854	13.745417	38.999896	34.645521
3	AlejandroSanz	40.377154	15.377462	52.337538	31.082154
4	alfredocasero1	36.664677	19.642258	48.530806	31.138871

- ✓ Definir previamente el número de cluster, es decir el número final de segmentos deseados

¿CÓMO DETERMINAR EL NÚMERO DE CLUSTER?



Elbow Point Example



- ✓ Iterar KMeans con un número cada vez mayor de cluster
- ✓ Identificar cuando la mejora del score sea mínima

Let's Code



<https://drive.google.com/file/d/16EOFMzByzl5GapnV30Cv7BAxleO81DRF/view?usp=sharing>