

APACHE SPARK

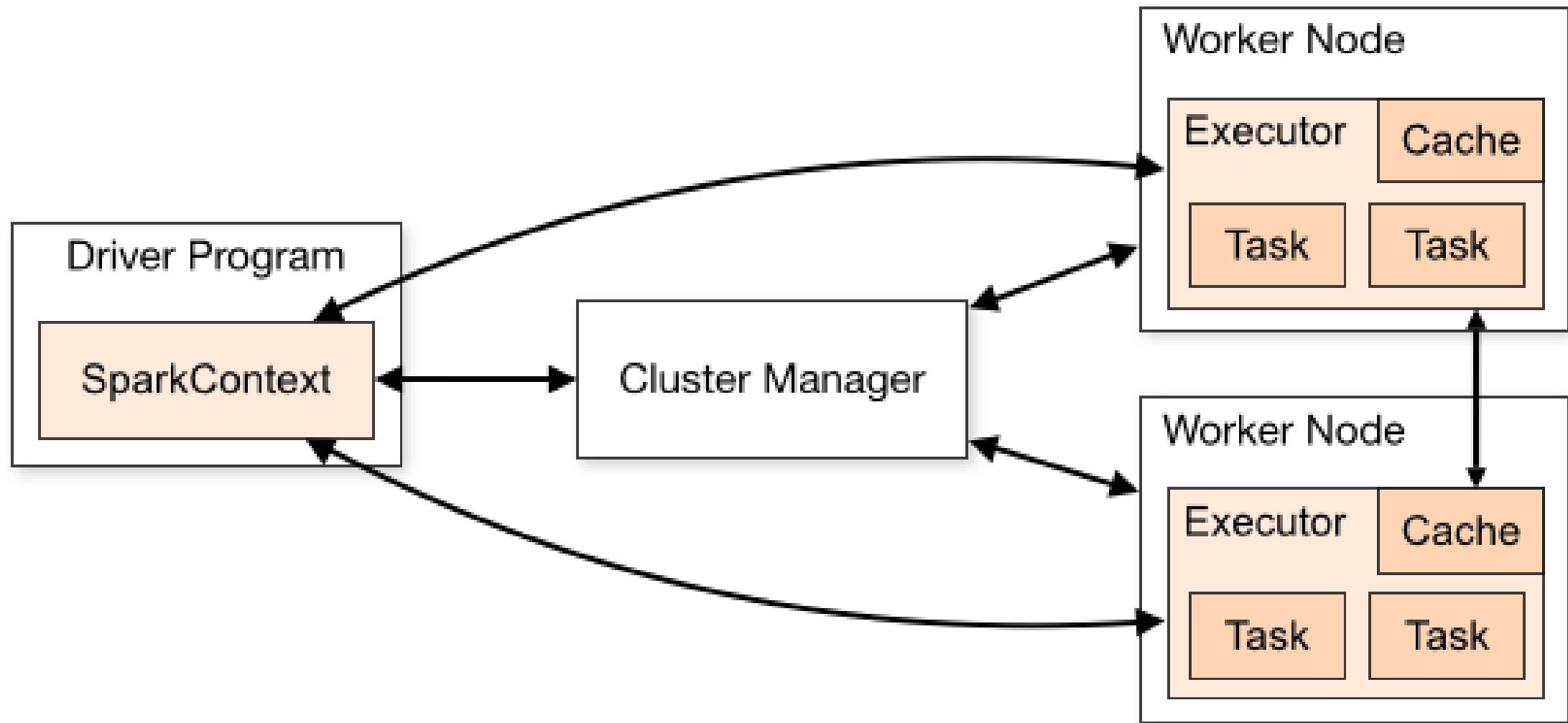


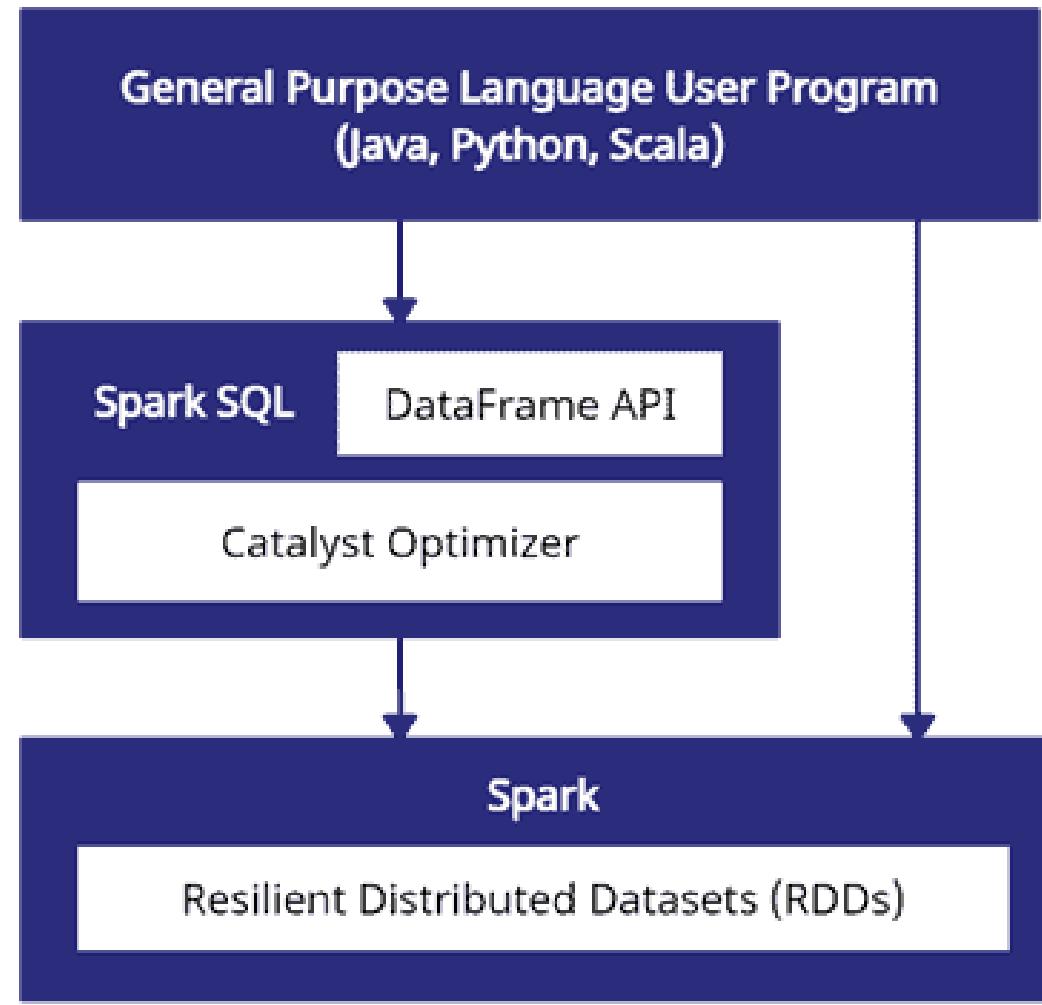
# ÍNDICE

- **Spark SQL**
- **MLlib**
  - Introducción a Machine Learning
  - Clasificaciones y Regresiones
  - Regresión: Regresión Lineal
  - Clasificación: Regresión Logística, Árboles de Decisión y Random Forest
  - Clustering: Kmeans
  - Sistemas de recomendación
- **Streaming**
- **GraphX**

# ■ RRD: Ciclo de Vida



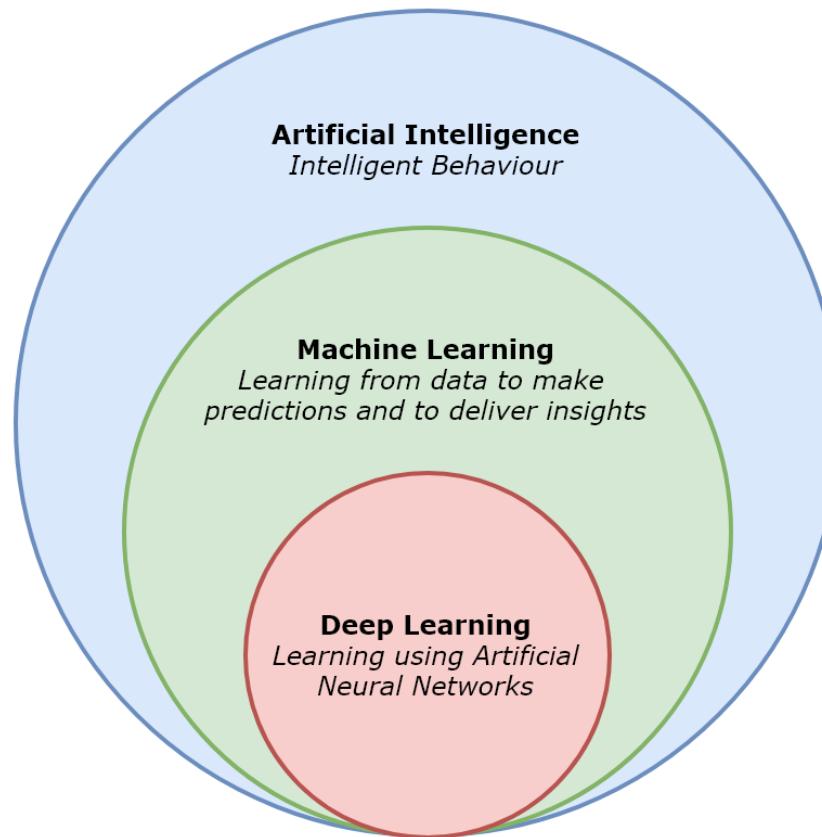




```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('style', 'border: 1px solid #ccc; border-radius: 5px; padding: 2px 10px; margin-bottom: 5px;');
55           tag.innerHTML = tag.name;
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK

# MACHINE LEARNING

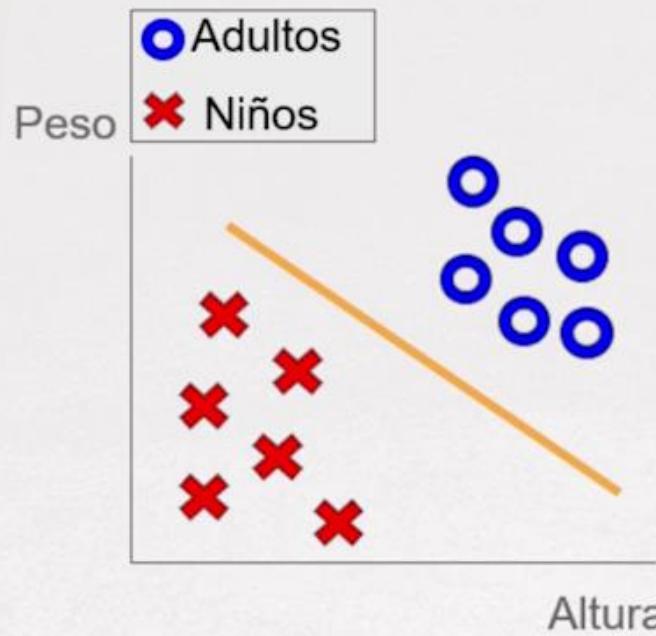


# TIPOS DE MACHINE LEARNING

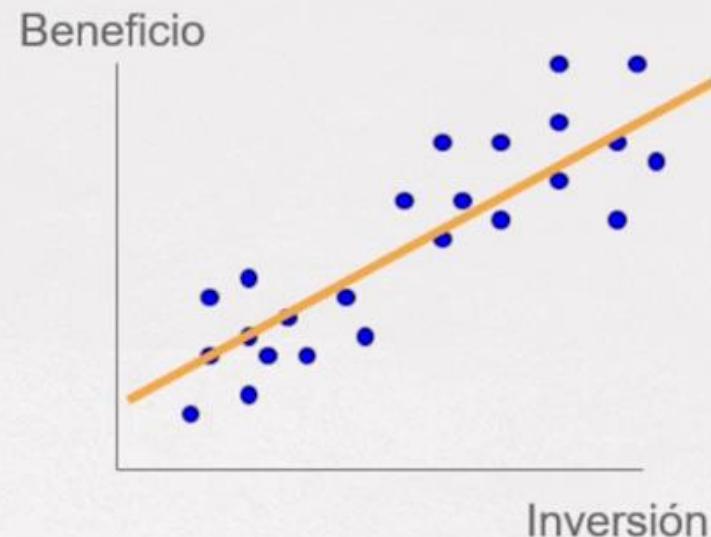


# Técnicas de Machine Learning

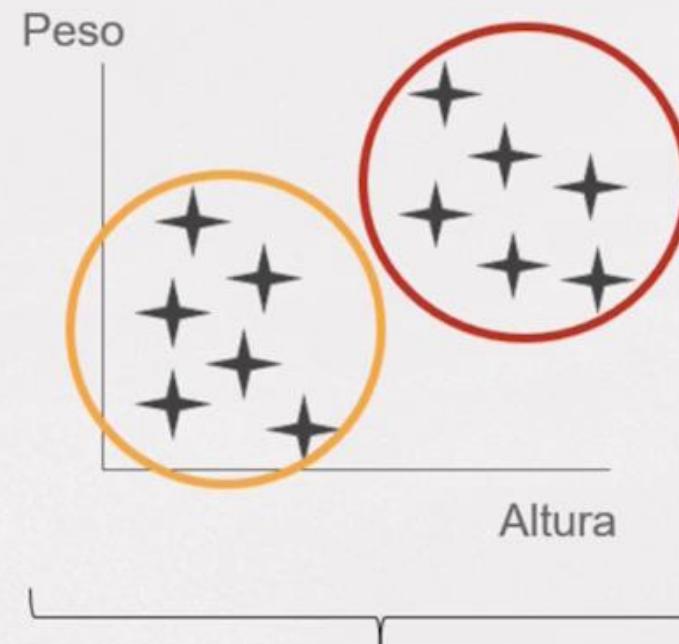
## Clasificación



## Regresión



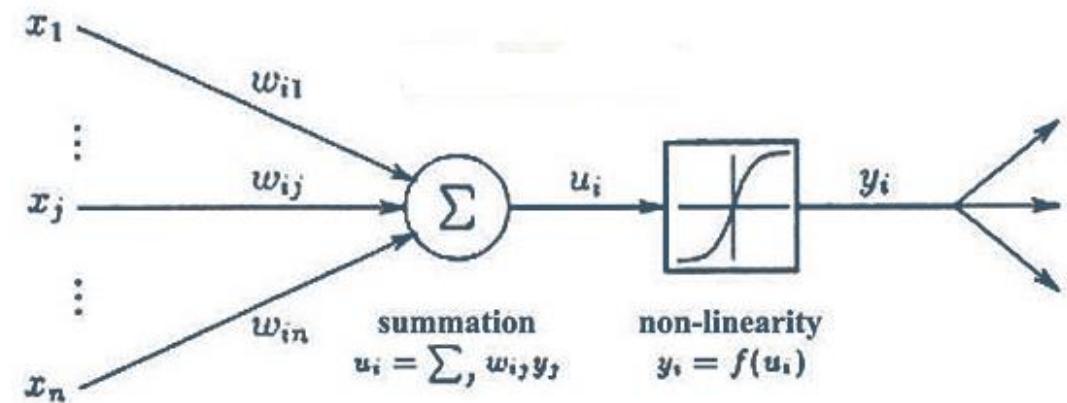
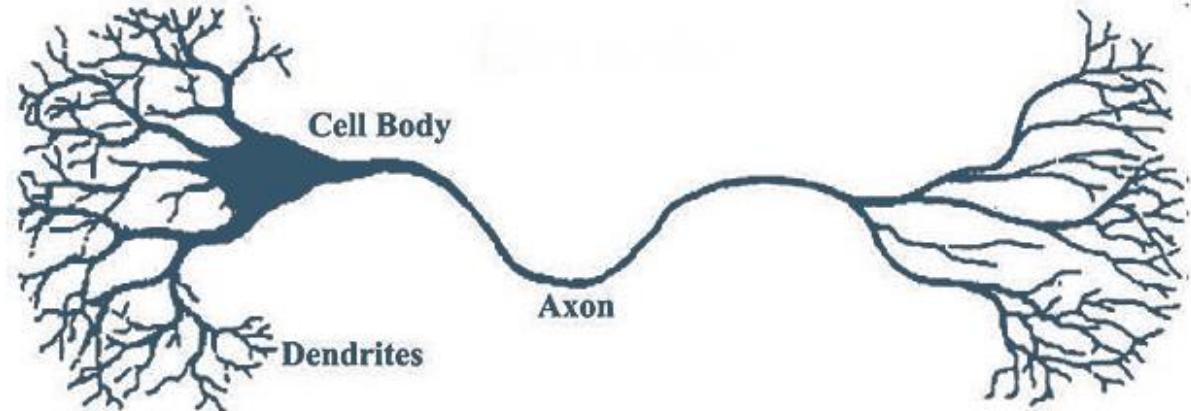
## Agrupación (*clustering*)



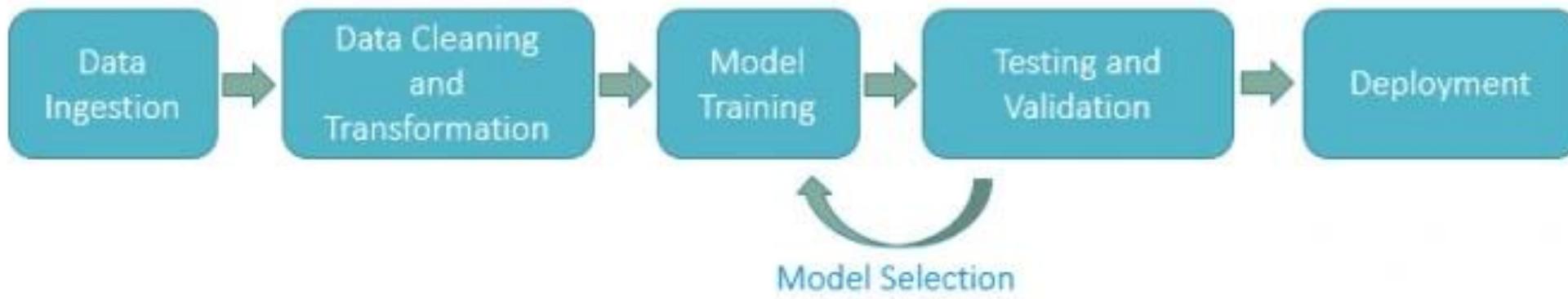
APRENDIZAJE SUPERVISADO

APRENDIZAJE NO SUPERVISADO

# REDES NEURONALES

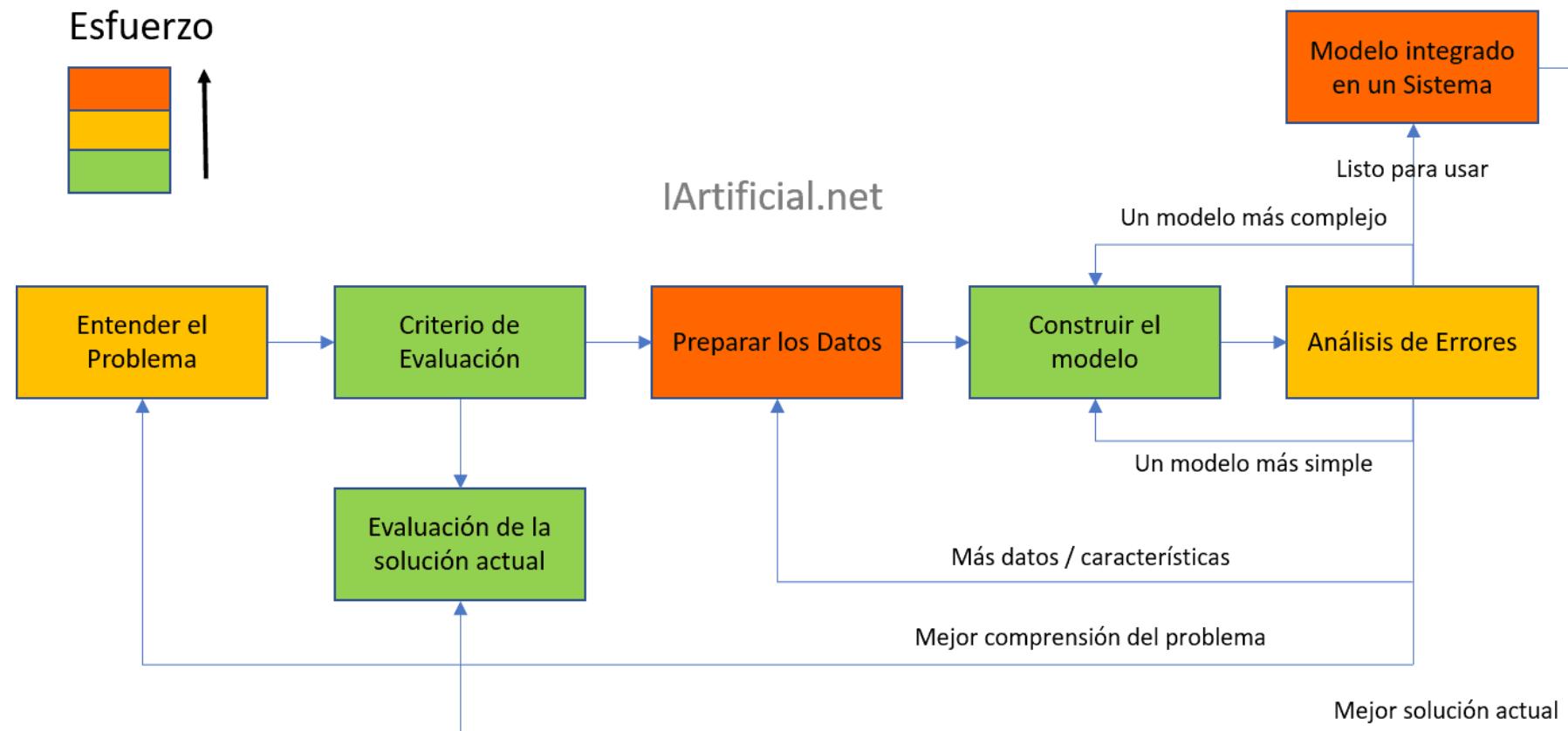


# APACHE Spark™ ML

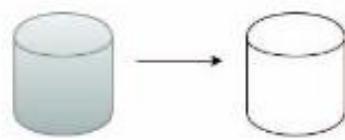


ML LIB

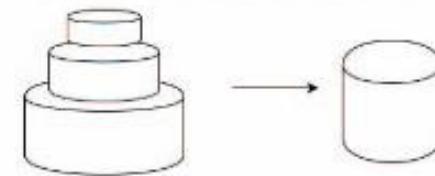
# PASOS DE UN PROYECTO DE ML



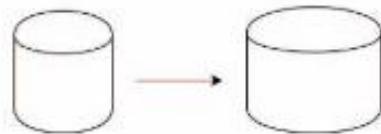
Limpieza de datos



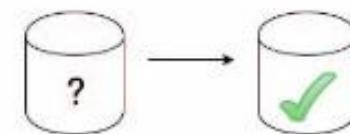
Normalización de datos



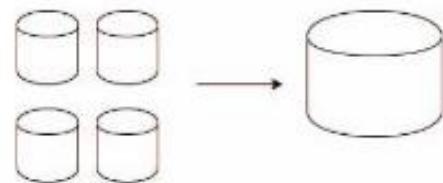
Transformación de datos



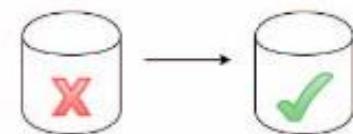
Imputación de valores perdidos



Integración de datos



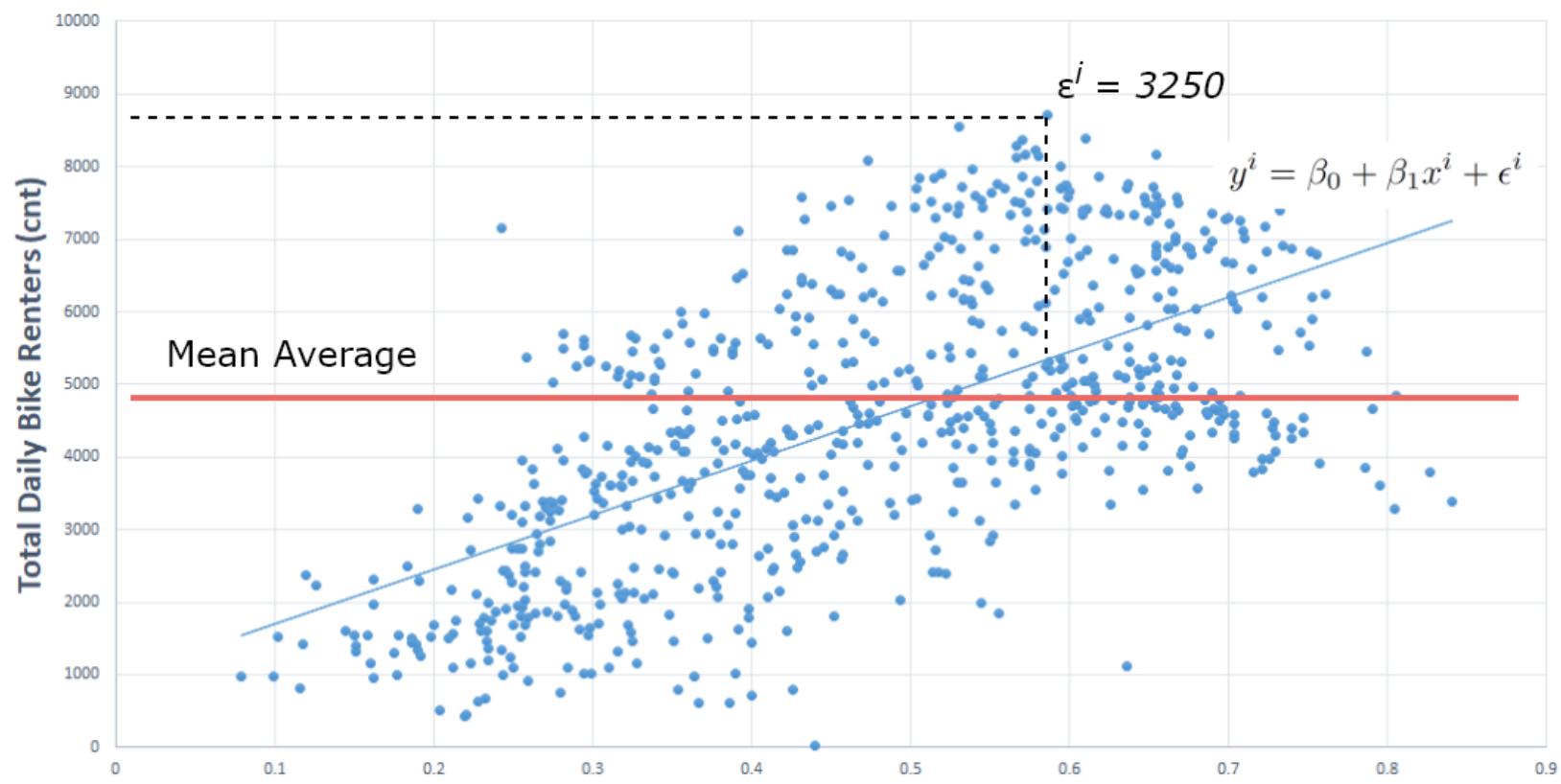
Identificación de ruido



# PROBLEMAS EN LOS DATOS

```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK



# REGRESIÓN LINEAL

$$y^i = \beta_0 + \beta_1 x_1^i + \beta_2 x_2^i + \dots + \beta_k x_k^i + \epsilon^i$$

$$SSE = (\epsilon_1)^2 + (\epsilon_2)^2 + \dots + (\epsilon_n)^2$$

$$R^2 = 1 - \frac{SSE}{SST}$$

# FÓRMULAS

# CORRELACIÓN



SOLÍA CREER QUE LA CORRELACIÓN IMPLICA CAUSALIDAD.



LUEGO DI UNA ASIGNATURA DE ESTADÍSTICA Y DEJÉ DE CREERLO.



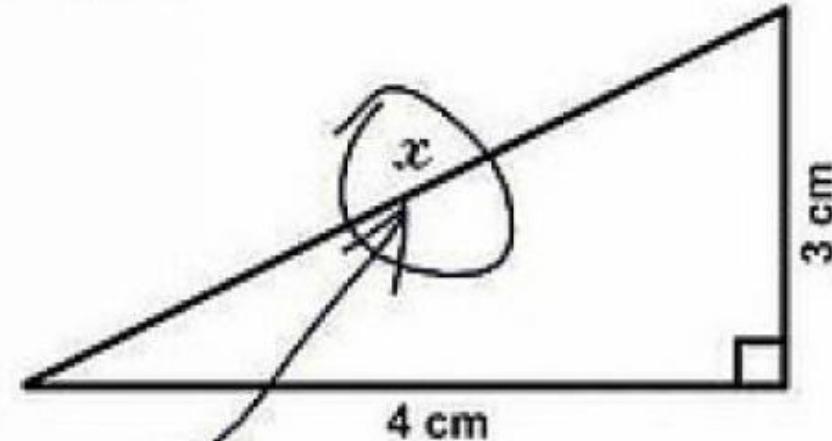
PARECE QUE ESA CLASE TE AYUDÓ.



BUENO, QUIZÁ.

# PRINCIPIO DE PARSIMONIA

3. Find x.

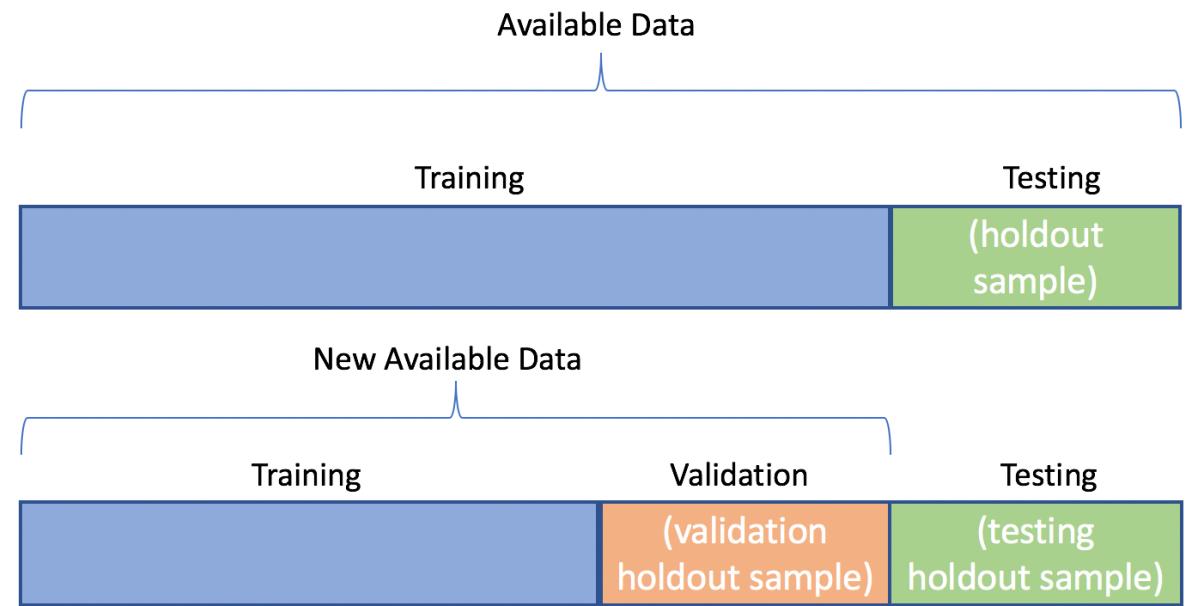


*Here it is*

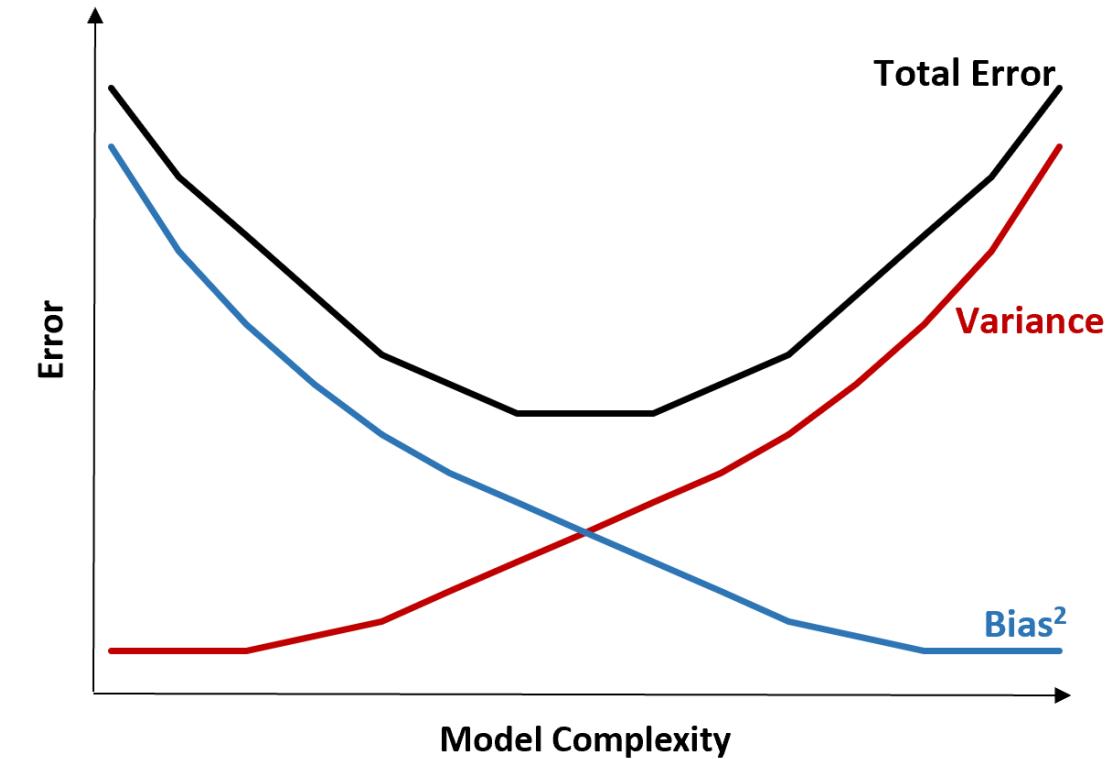
desmotivaciones.es

Principio de parsimonia  
la solución más simple suele ser la mejor.

# TRAIN Y TEST

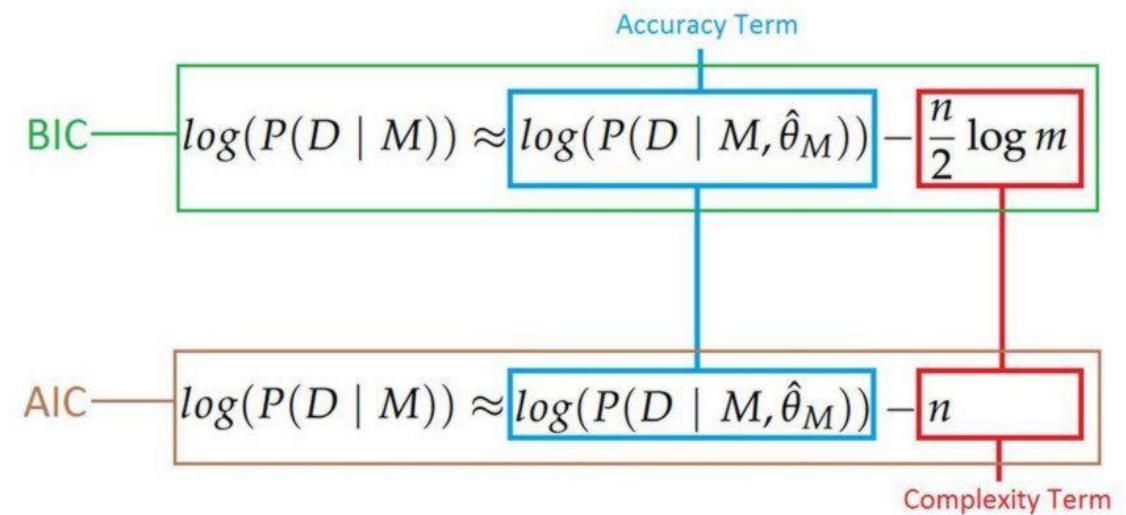


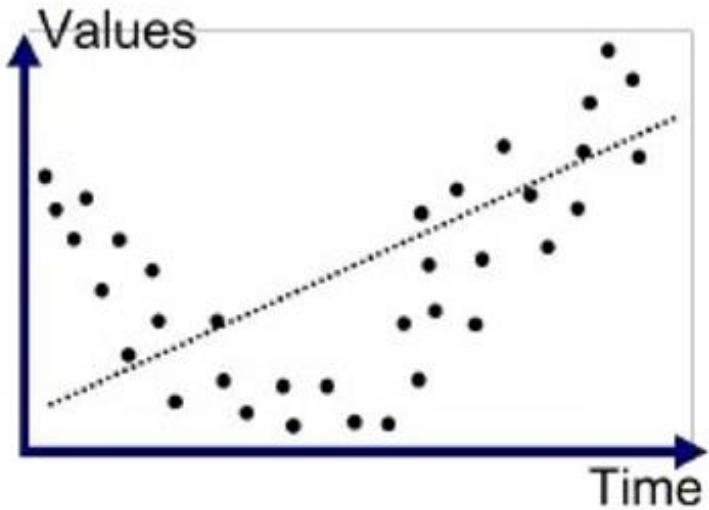
# TRADE OFF SESGO Y VARIANZA



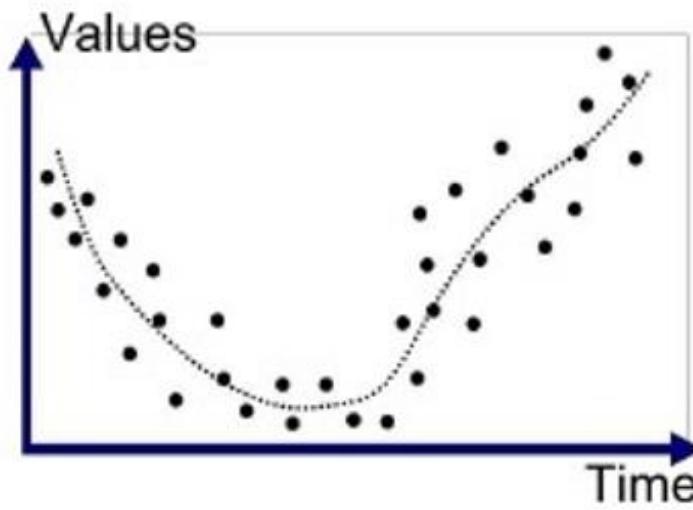
# AIC VS BIC

## AIC vs BIC

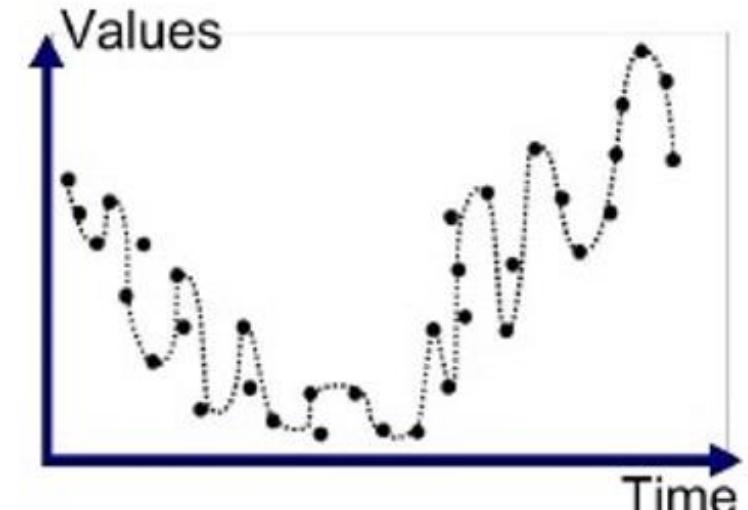




Underfitted

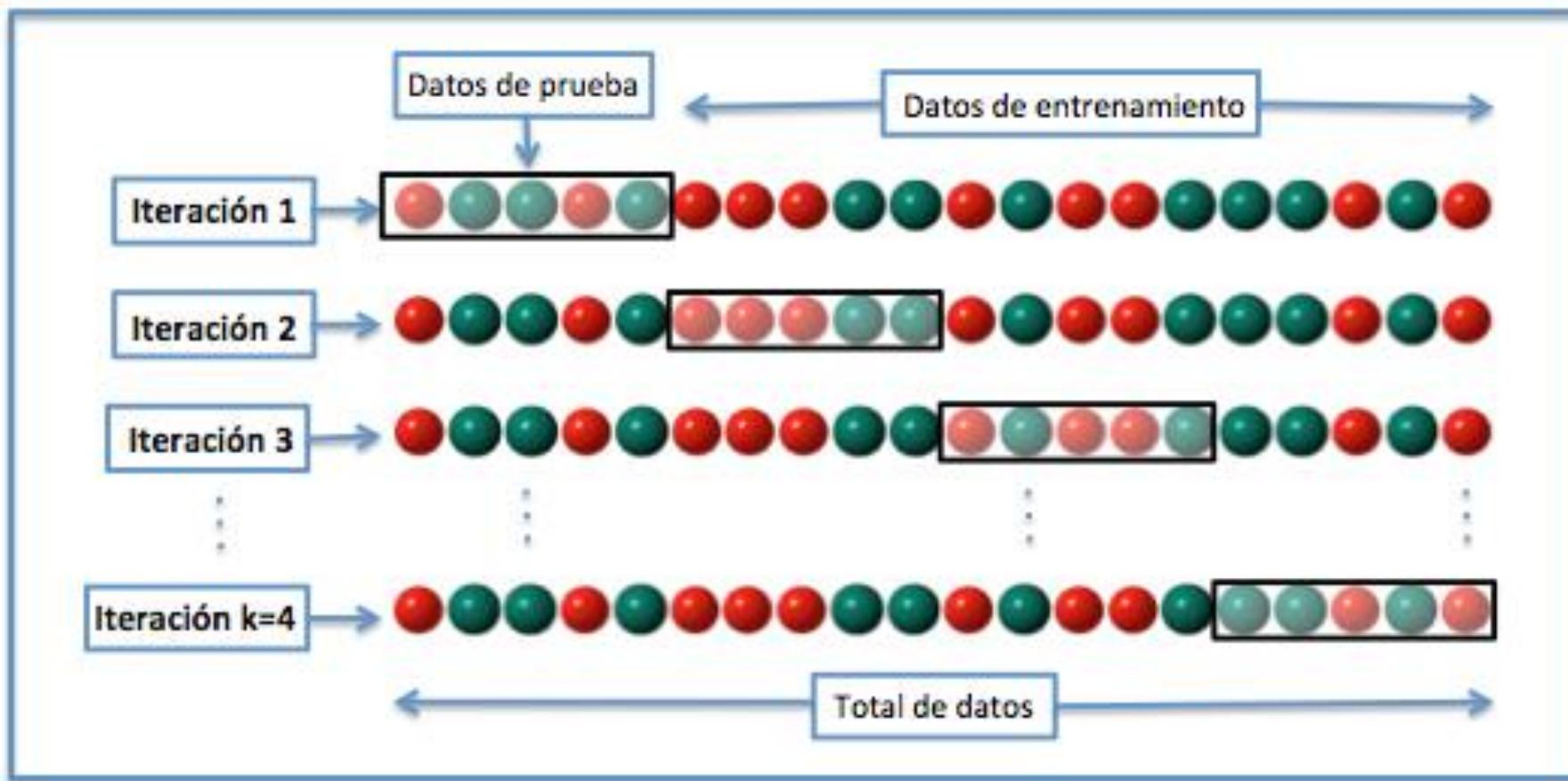


Good Fit/Robust



Overfitted

LASSO, RIDGE Y ELASTIC NET



# VALIDACION CRUZADA

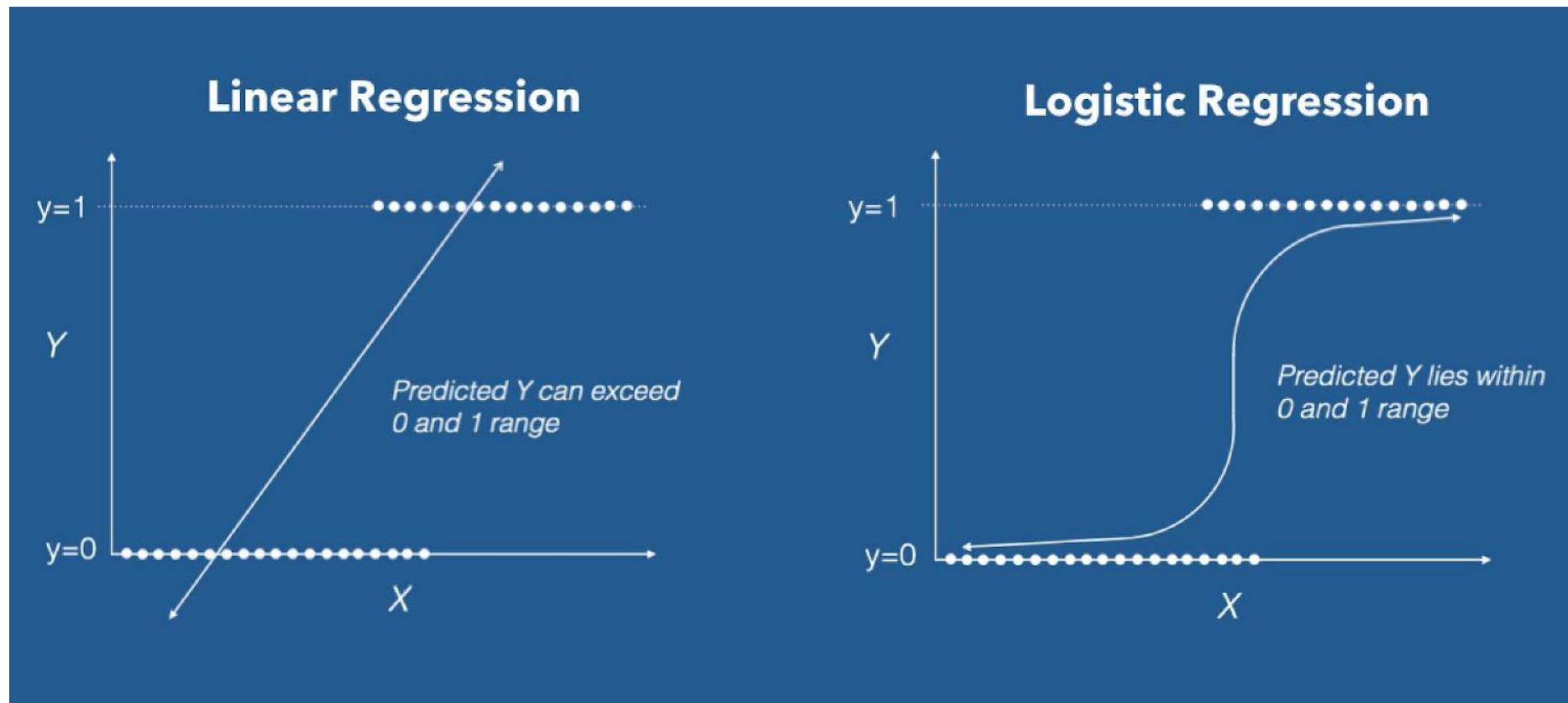
```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK

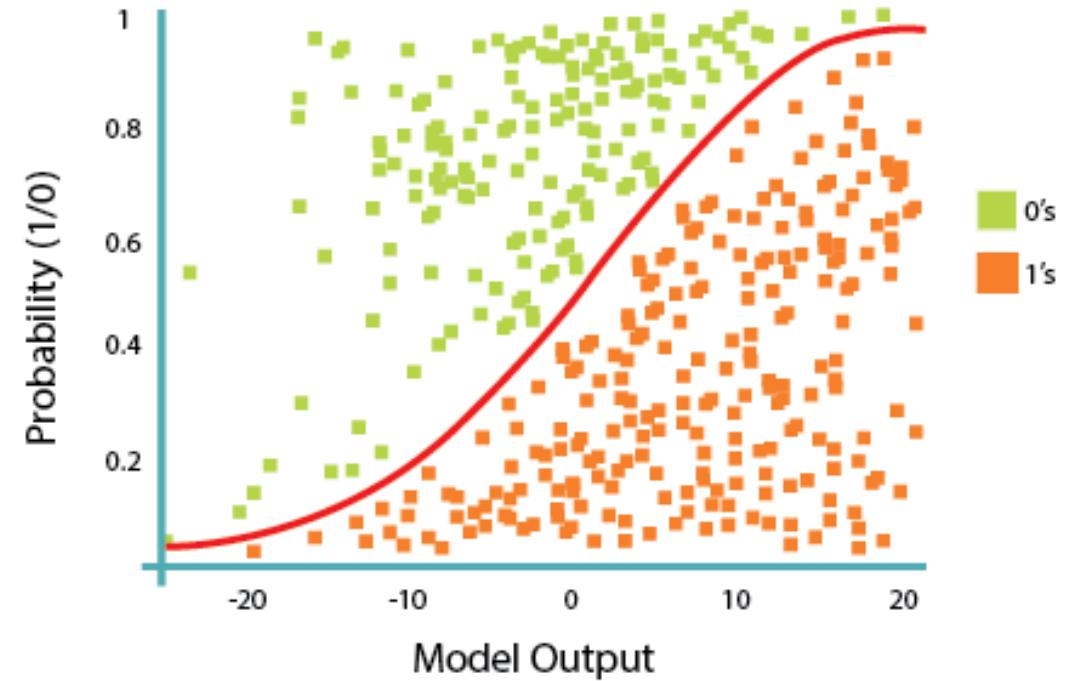
```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK

# REGRESIÓN LOGÍSTICA



# FUNCIÓN SIGMOIDE

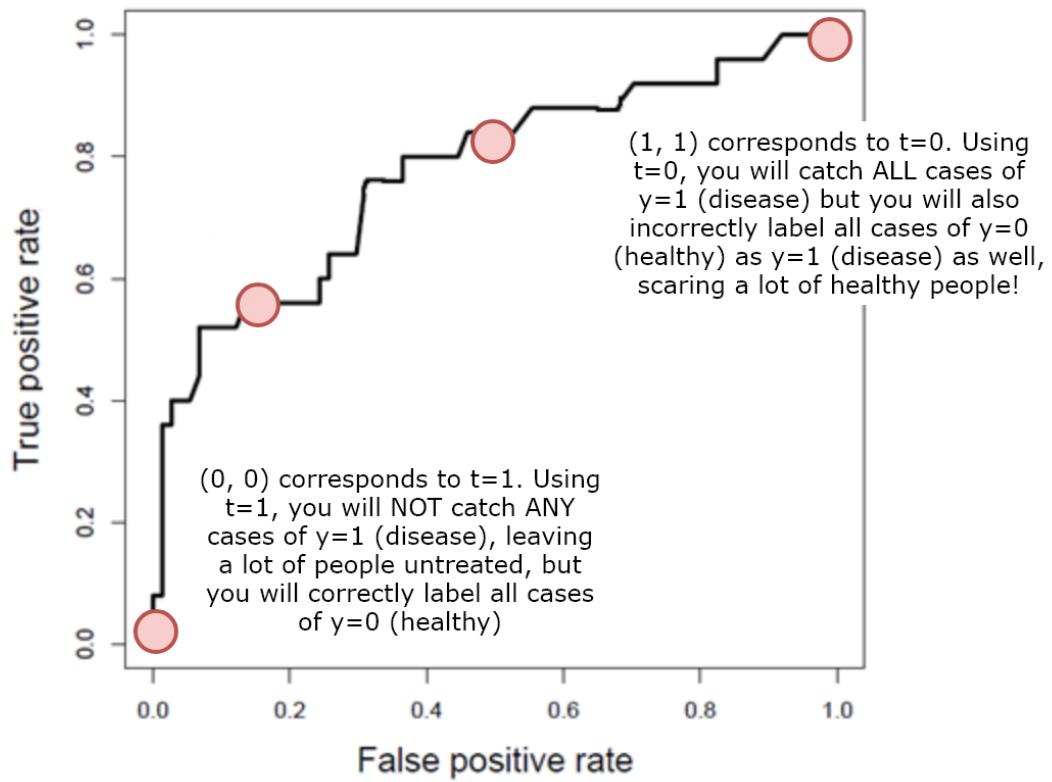


# THRESHOLD

		Reality	
		True	False
Measured or Perceived	True	Correct 😊	Type 1 error False Positive
	False	Type 2 error False Negative	Correct 😊

Matriz de confusión		Estimado por el modelo				
		Negativo (N)	Positivo (P)			
Real	Negativo	a: (TN)	b: (FP)	<b>Precisión</b> ("precision") Porcentaje predicciones positivas correctas:	$d/(b+d)$	
	Positivo	c: (FN)	d: (TP)			
		<b>Sensibilidad, exhaustividad ("Recall")</b> Porcentaje casos positivos detectados	<b>Especificidad (Specificity)</b> Porcentaje casos negativos detectados	<b>Exactitud ("accuracy")</b> Porcentaje de predicciones correctas <i>(No sirve en datasets poco equilibrados)</i>		
		$d/(d+c)$	$a/(a+b)$	$(a+d)/(a+b+c+d)$		

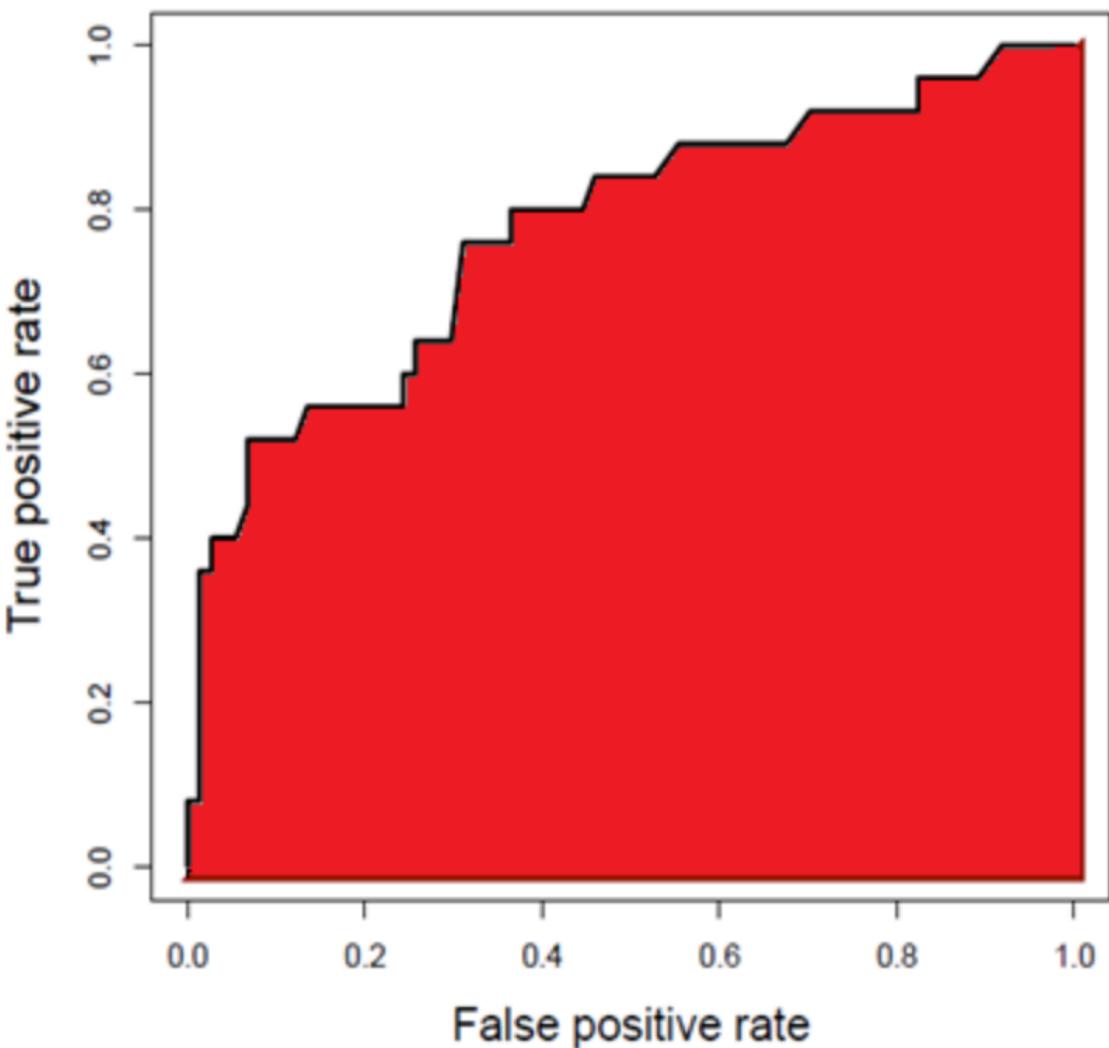
### Receiver Operator Characteristic Curve



CURVA ROC

AUC

Receiver Operator Characteristic Curve



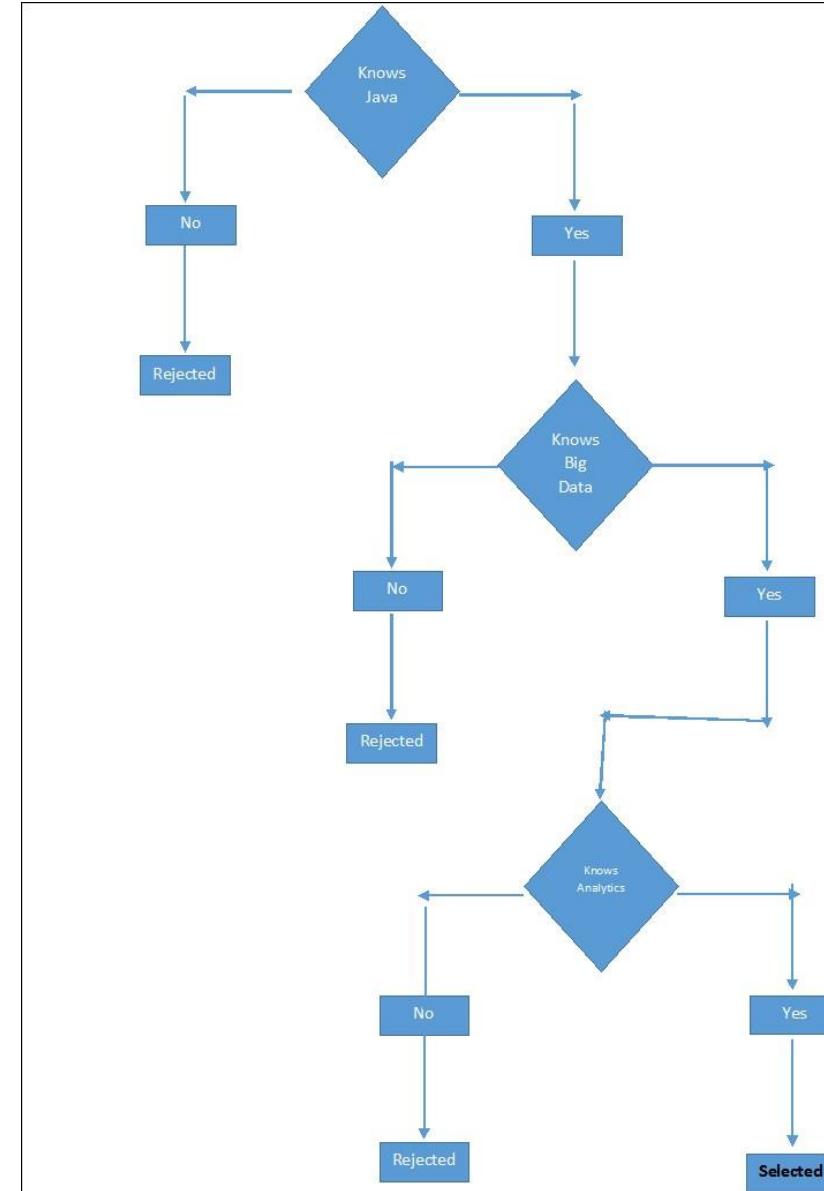
```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK

```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

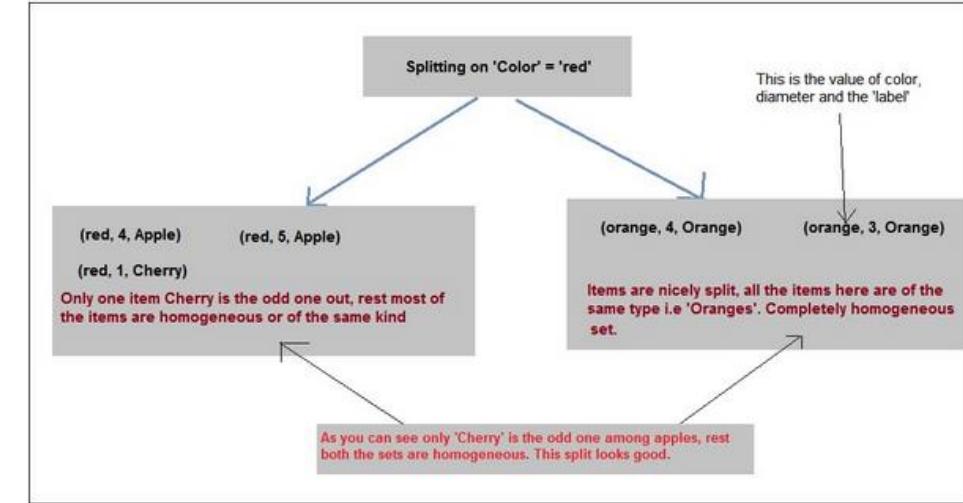
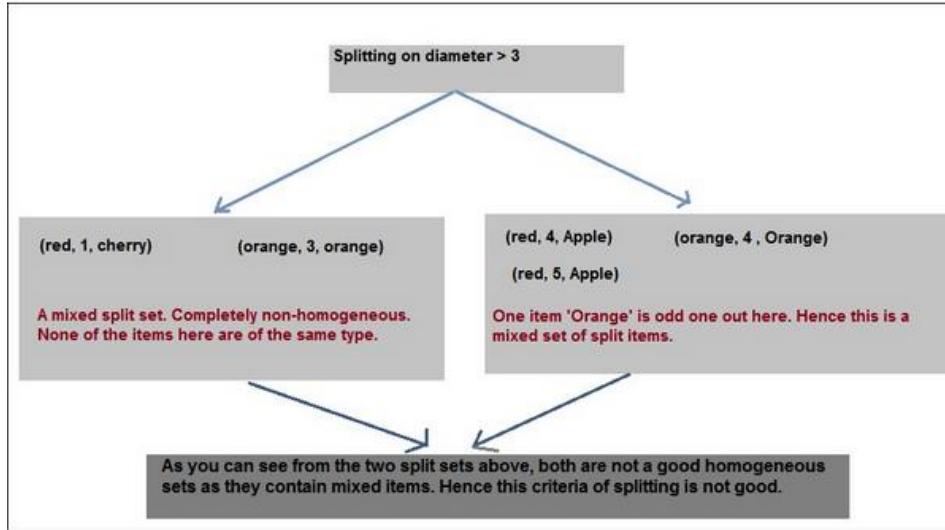
# NOTEBOOK

# ARBOLES DE DECISIÓN



# EJEMPLO DE ARBOL DE DECISIÓN

Color	Diámetro	Fruta
Rojo	4	Manzana
Naranja	4	Naranja
Rojo	1	Cereza
Rojo	5	Manzana
Naranja	3	Naranja



# COMPARANDO DIVISIONES

$$Entropy(fruits) = - \left[ \left( \frac{3}{8} \right) \log_2 \frac{3}{8} + \left( \frac{5}{8} \right) \log_2 \frac{5}{8} \right] = 0.95$$

$$Gini\ Impurity = 1 - \left[ \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right] = 0.48$$

## INDICE DE GINI Y ENTROPIA

# VENTAJAS Y DESVENTAJAS

- ✓ Son muy simples de construir y usar. Son esencialmente un con-junto de if...else que conducen a un resultado concluyente.
- ✓ La entrada puede ser de cualquier tipo, ya sea numérica o string.
- ✓ Desde la perspectiva de los grandes datos y la informática distribuida, es más fácil crear un modelo de árbol de decisiones que se pueda distribuir en un grupo de máquinas. Por lo tanto, puede ejecutarse en paralelo y puede ser muy rápido.
- ❖ Los árboles de decisión sufren el problema del sobreajuste. Debido al sobreajuste, los se desempeñan extremadamente bien con los datos de entrenamiento, pero se desempeñan mal con cualquier dato nuevo que no sea parte del conjunto de entrenamiento.
- ❖ Si hay demasiadas reglas de decisión, el modelo pronto puede volverse bastante complejo. Dado que estamos tratando con grandes datos, este problema es más común.

# KINNECT

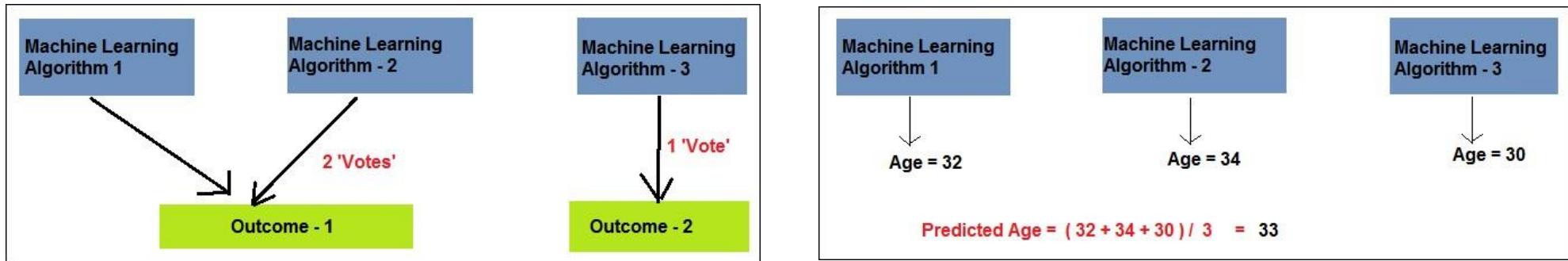


```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK

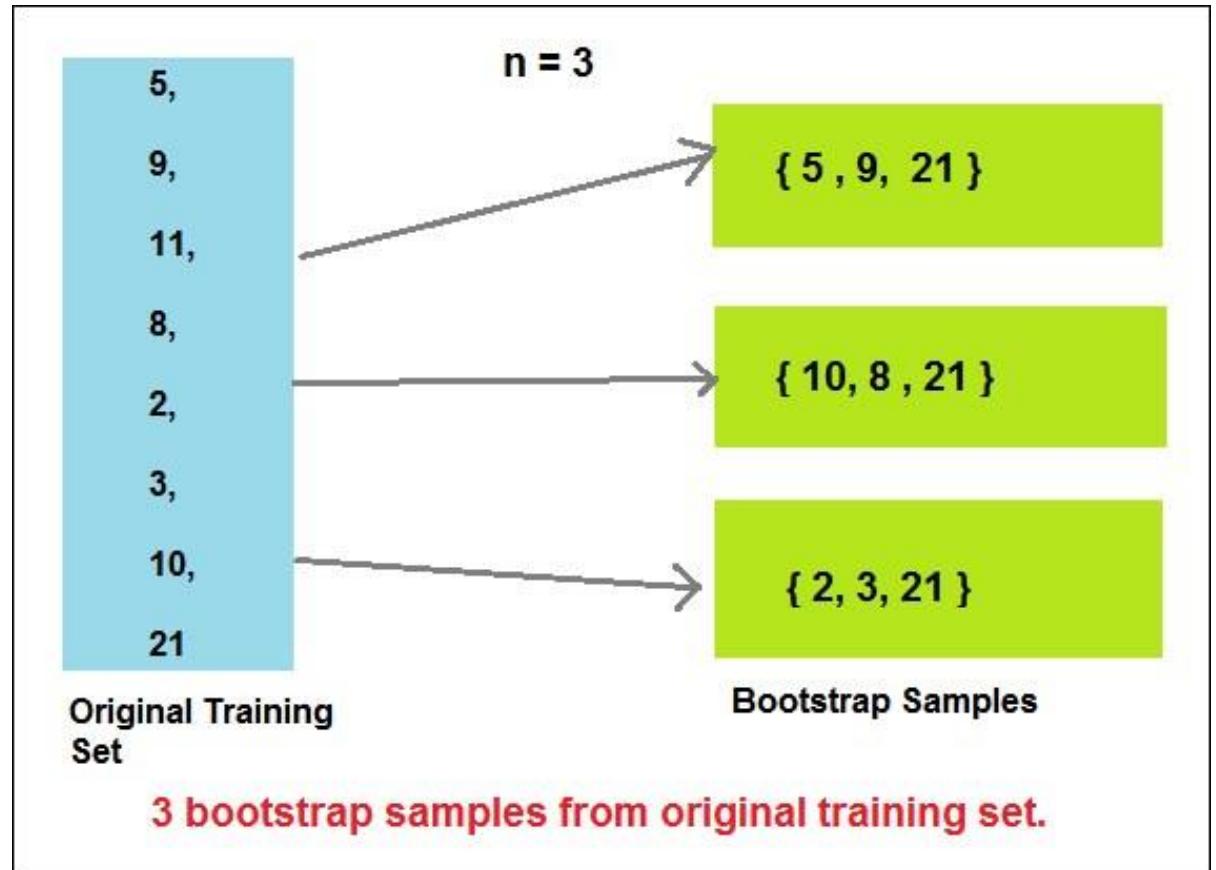
ENSEMBLING



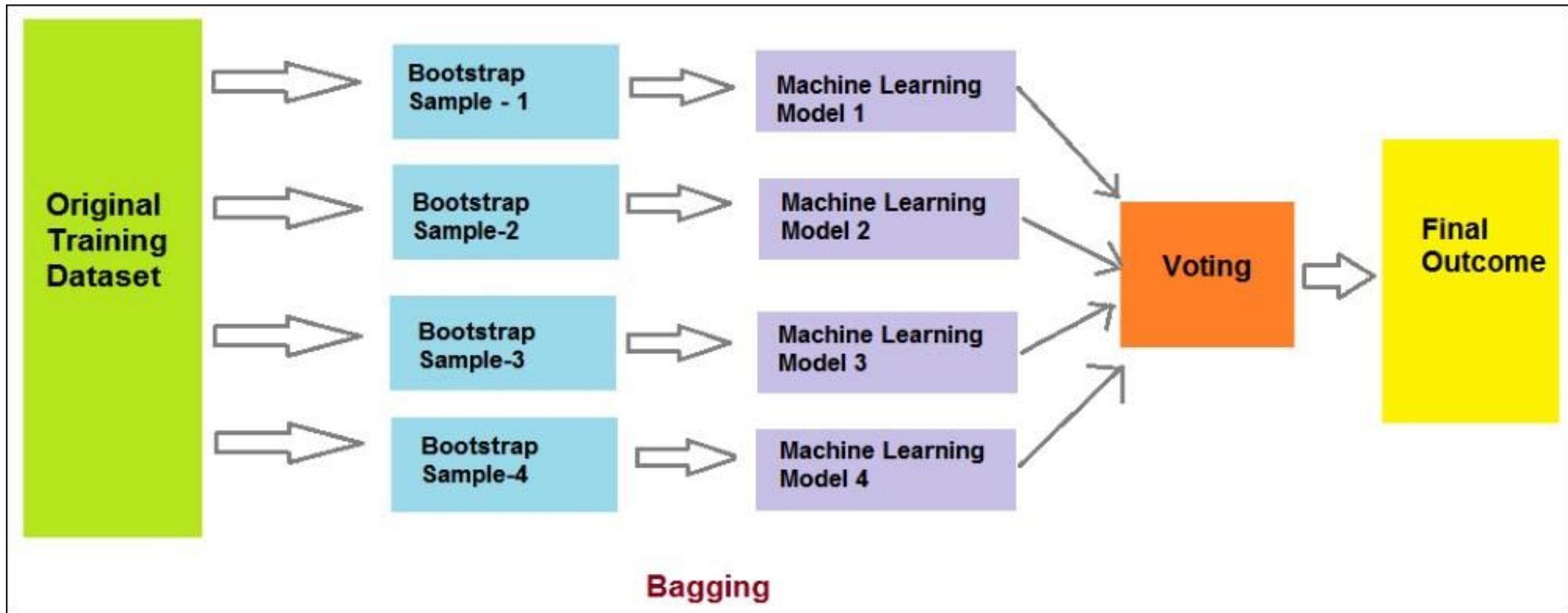


# ENFOQUES

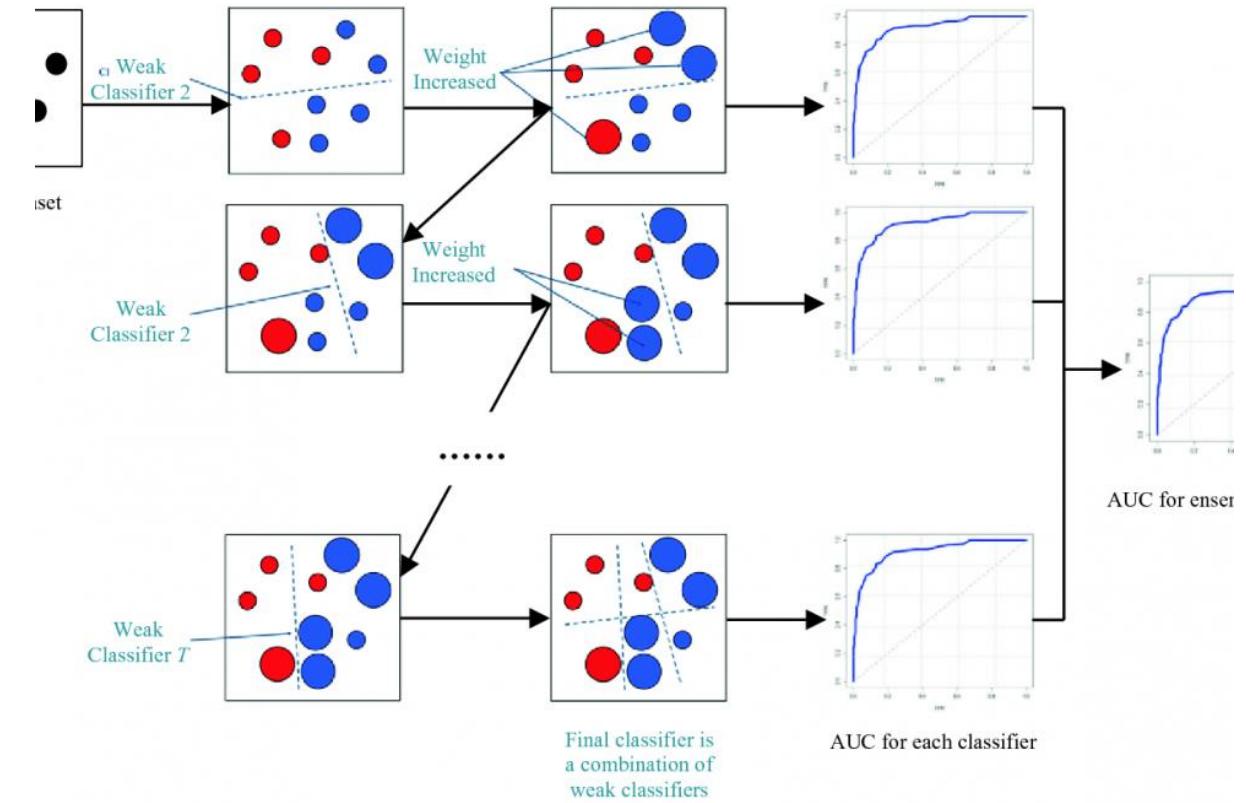
# BOOTSTRAPING



# BAGGING

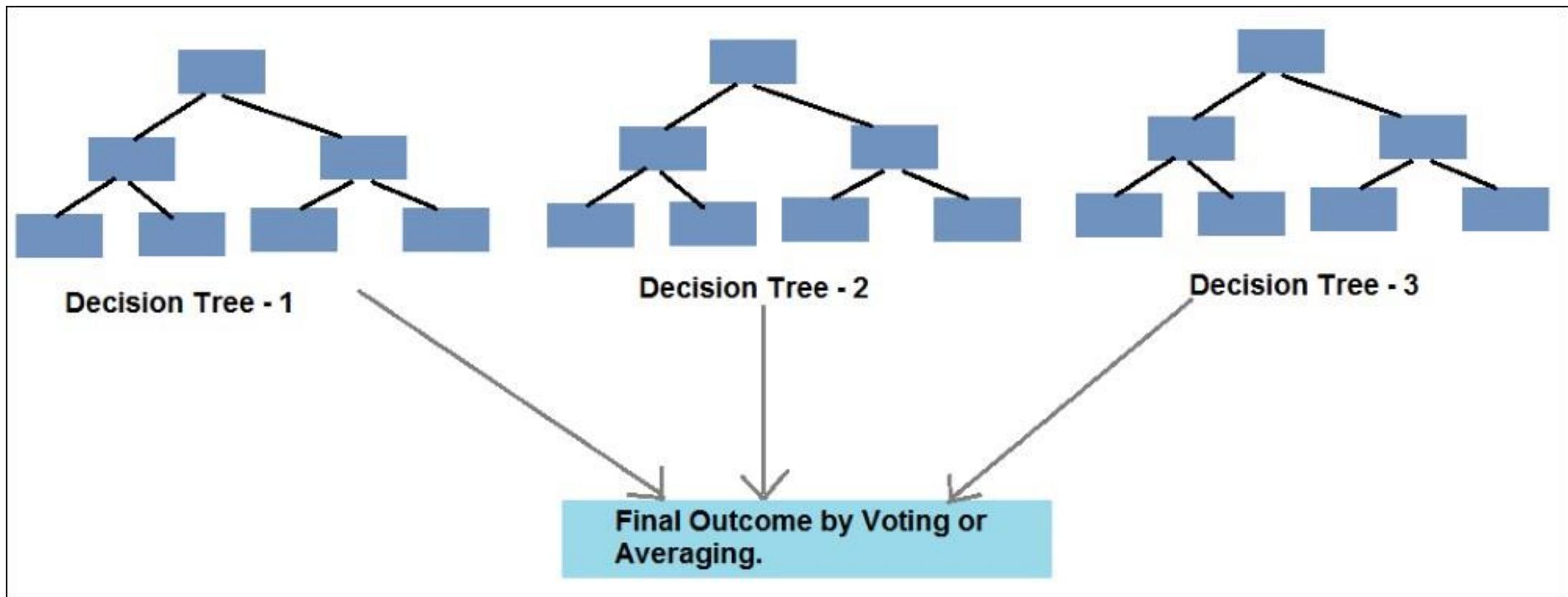


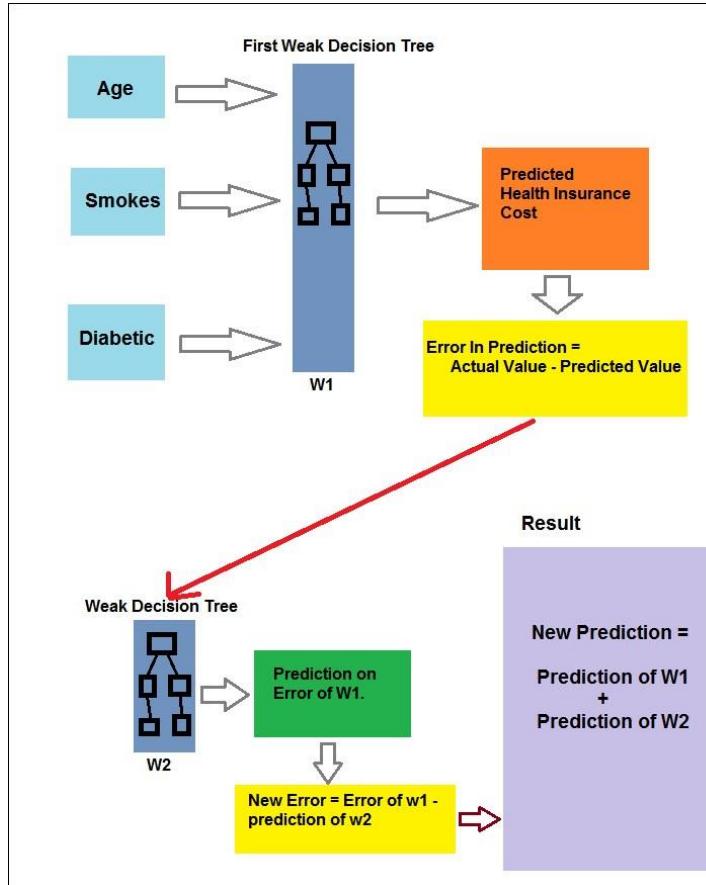
# BOOSTING



# VENTAJAS Y DESVENTAJAS

# RANDOM FOREST

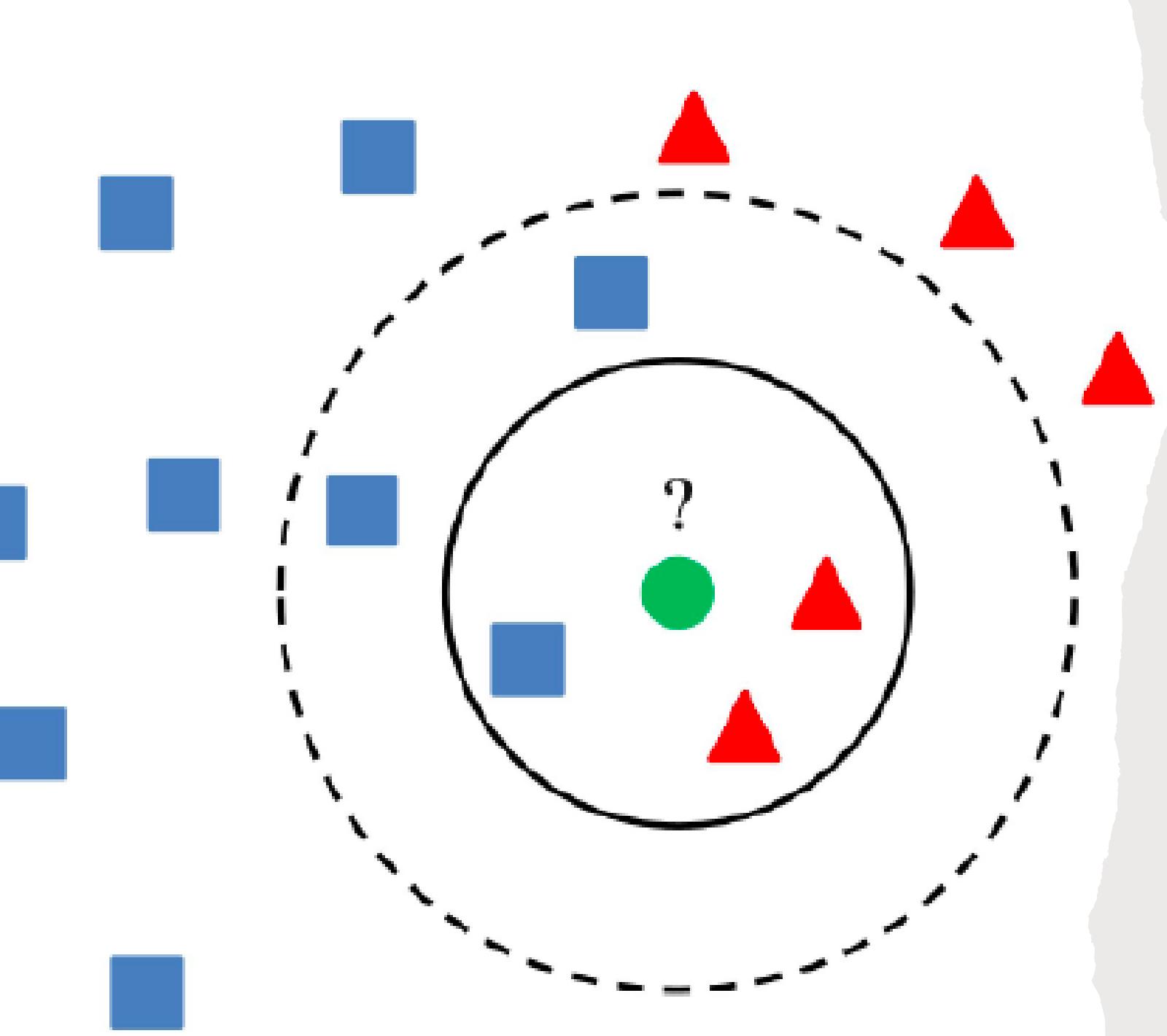




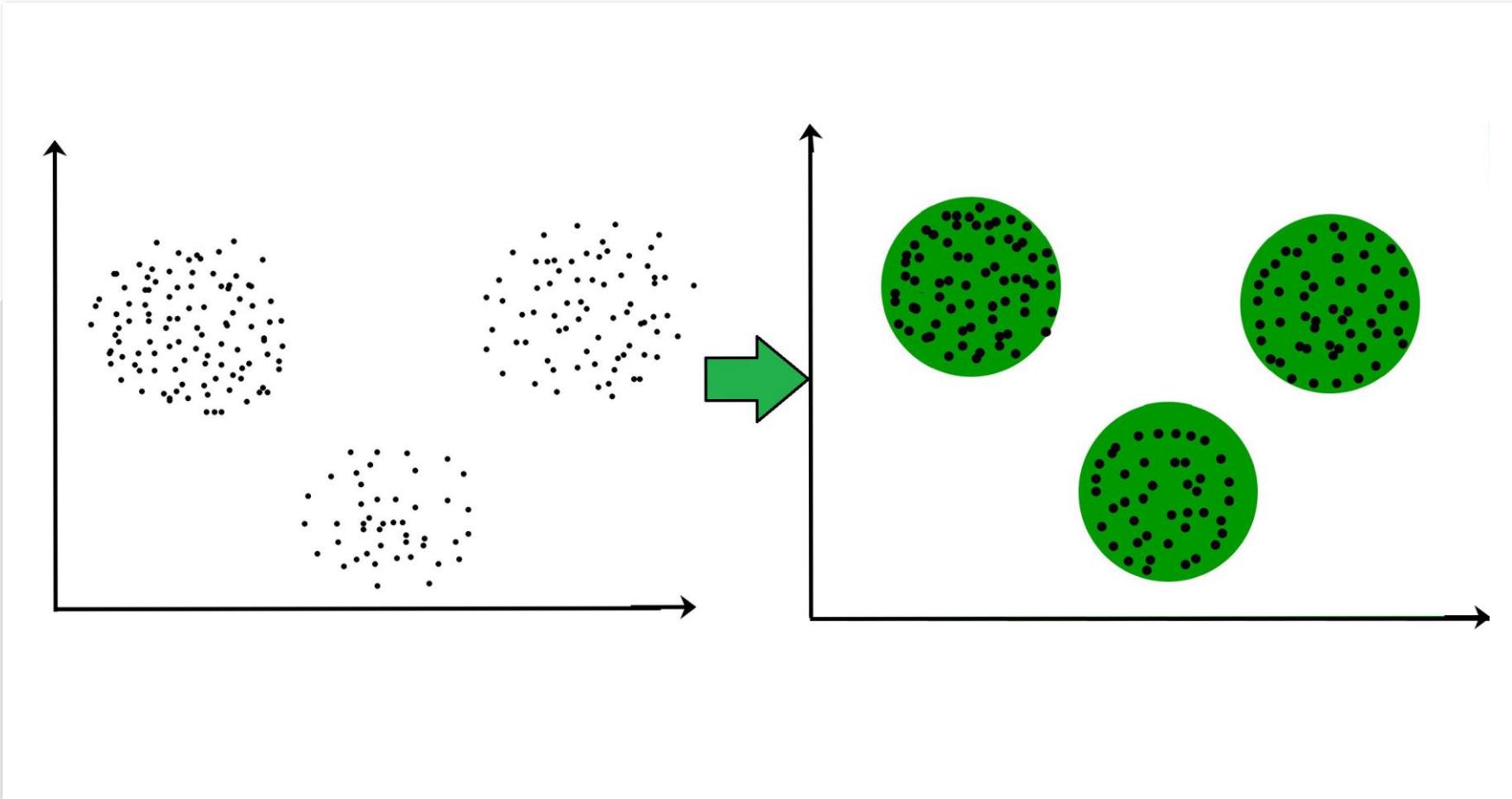
# GBT

```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK



KNN



CLUSTERING

All Departments ▾ car

1-16 of 6,153 results for "car"

Show results for

**Automotive >**

- Automotive Interior Accessories
- Cleaners
- Passenger Car Tires
- Bumper Stickers, Decals & Magnets
- Seat Cover Accessories

**Toys & Games >**

- Ride-On Toys
- Kids' Electronics
- Hobbies
- Children's Die-Cast Vehicles
- Hobby RC Cars

**Apps & Games >**

- Simulation Games
- Racing Games
- Action Games

**Electronics >**

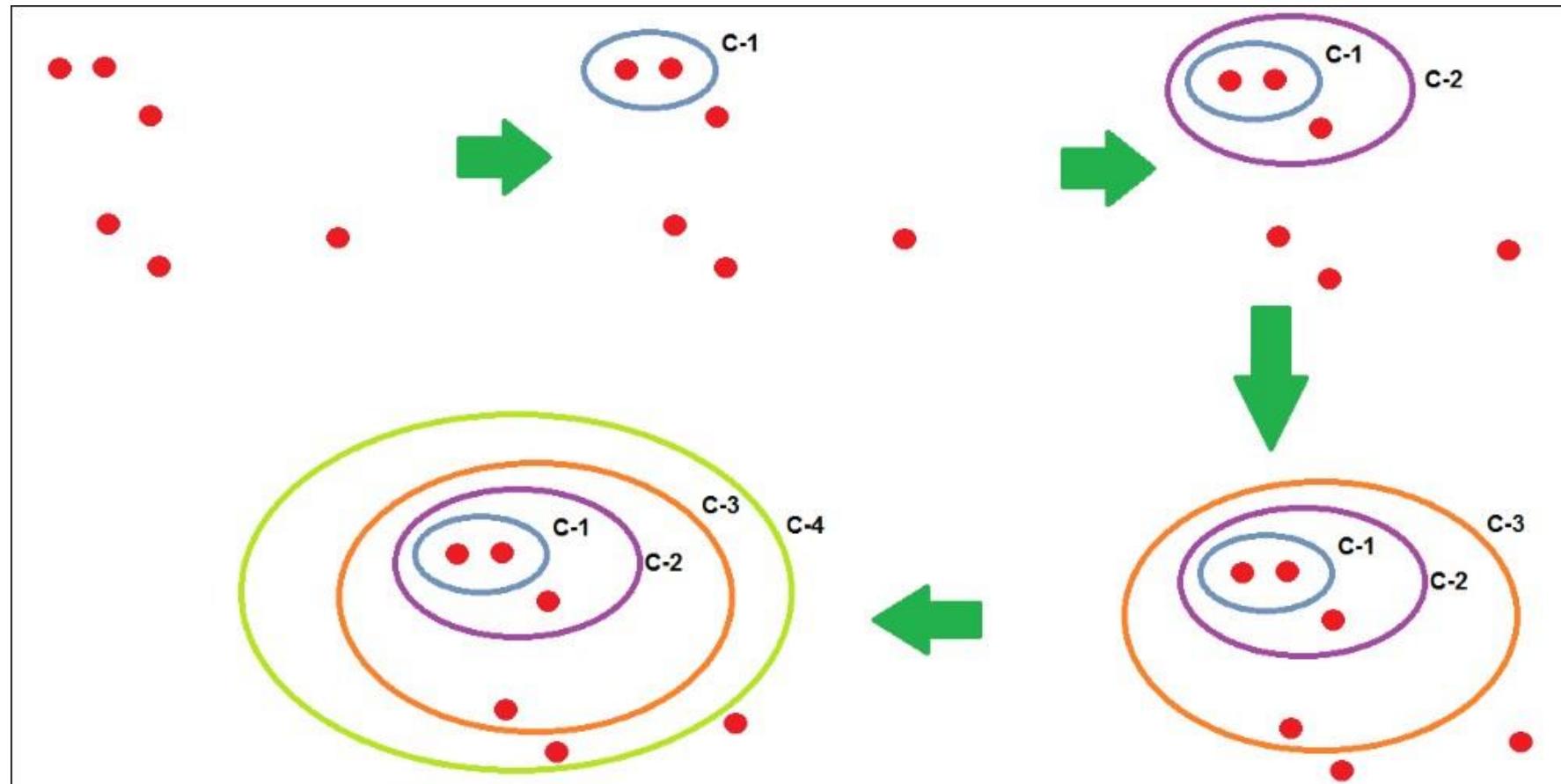
 Save Now on Swiffer

  
Haktoys HAK139 UTV SSV ATV 1:12 Scale RC Car with Lights

**Searched for 'CAR'**

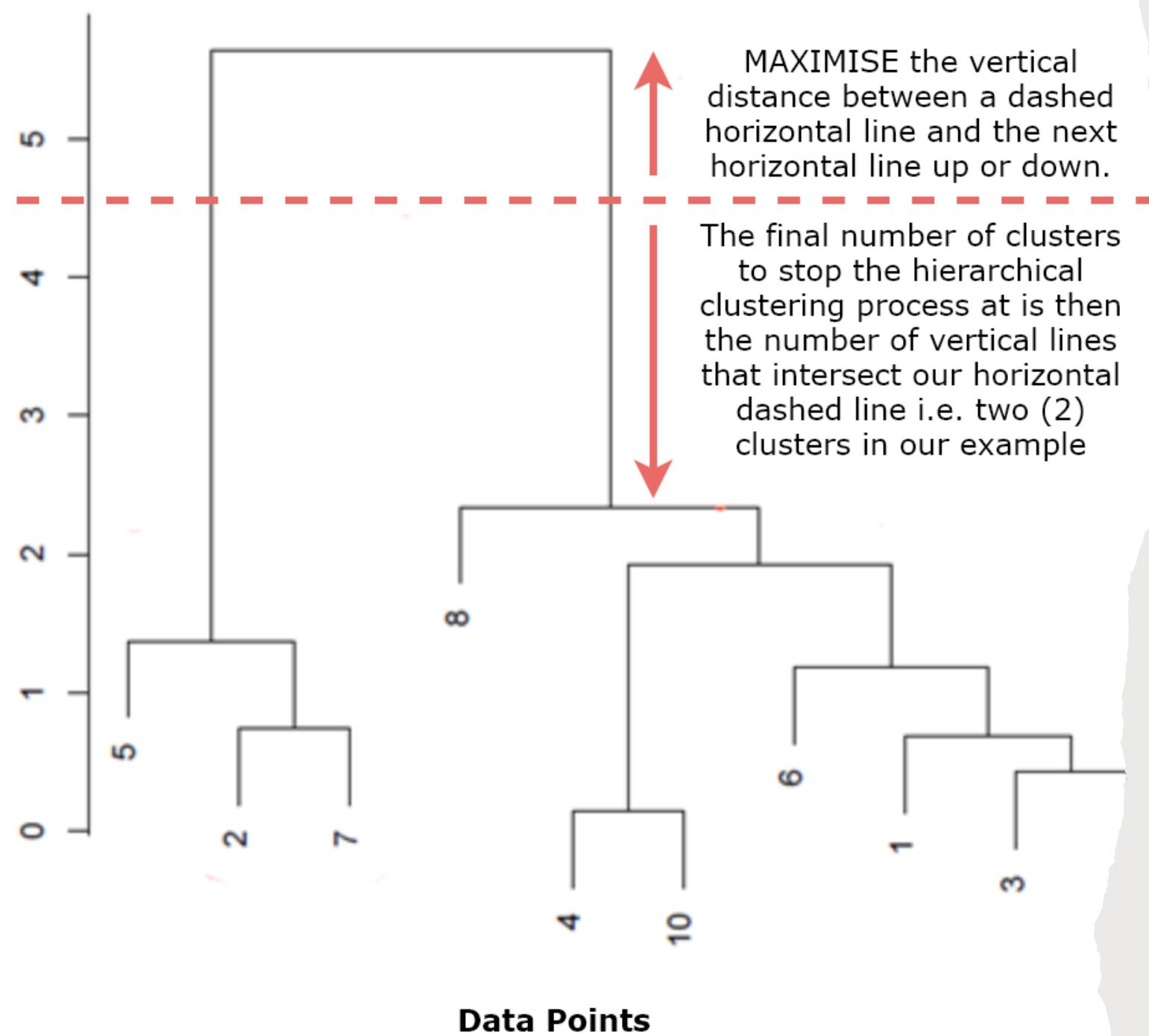
**Clustering on the word 'CAR' to figure out individual categories.**

# CASOS DE USO

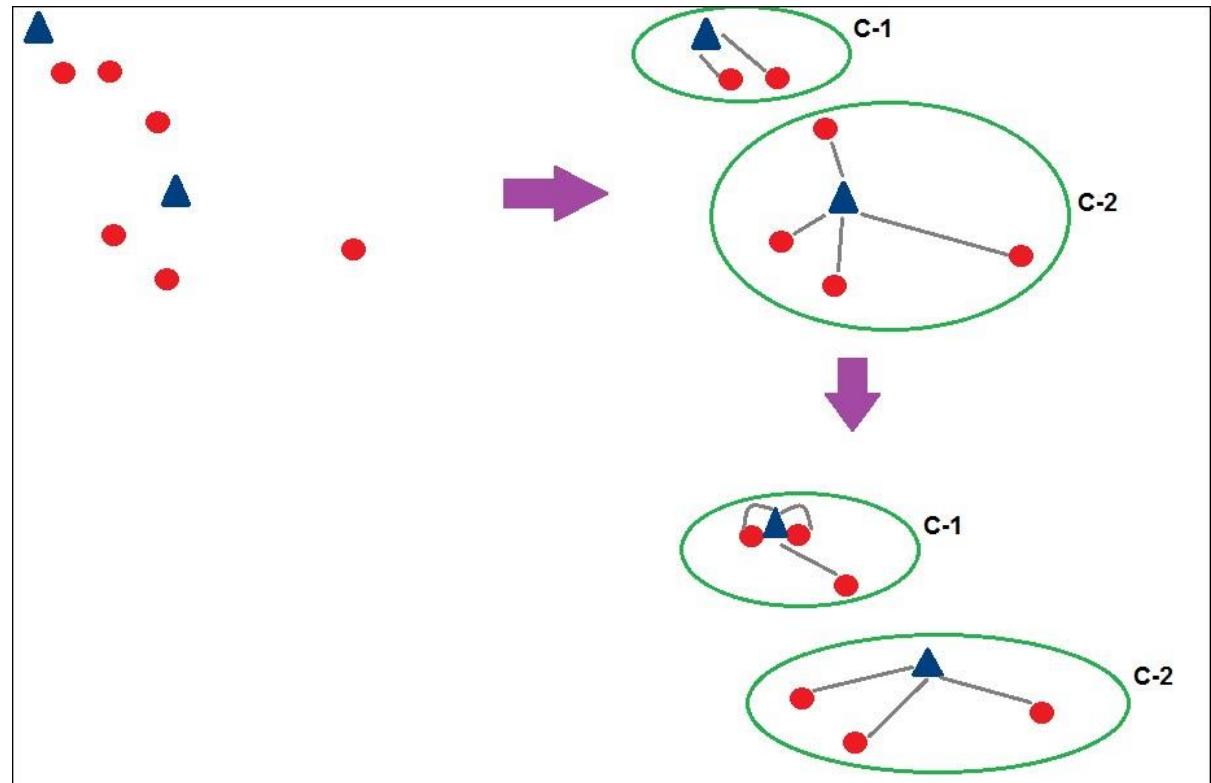


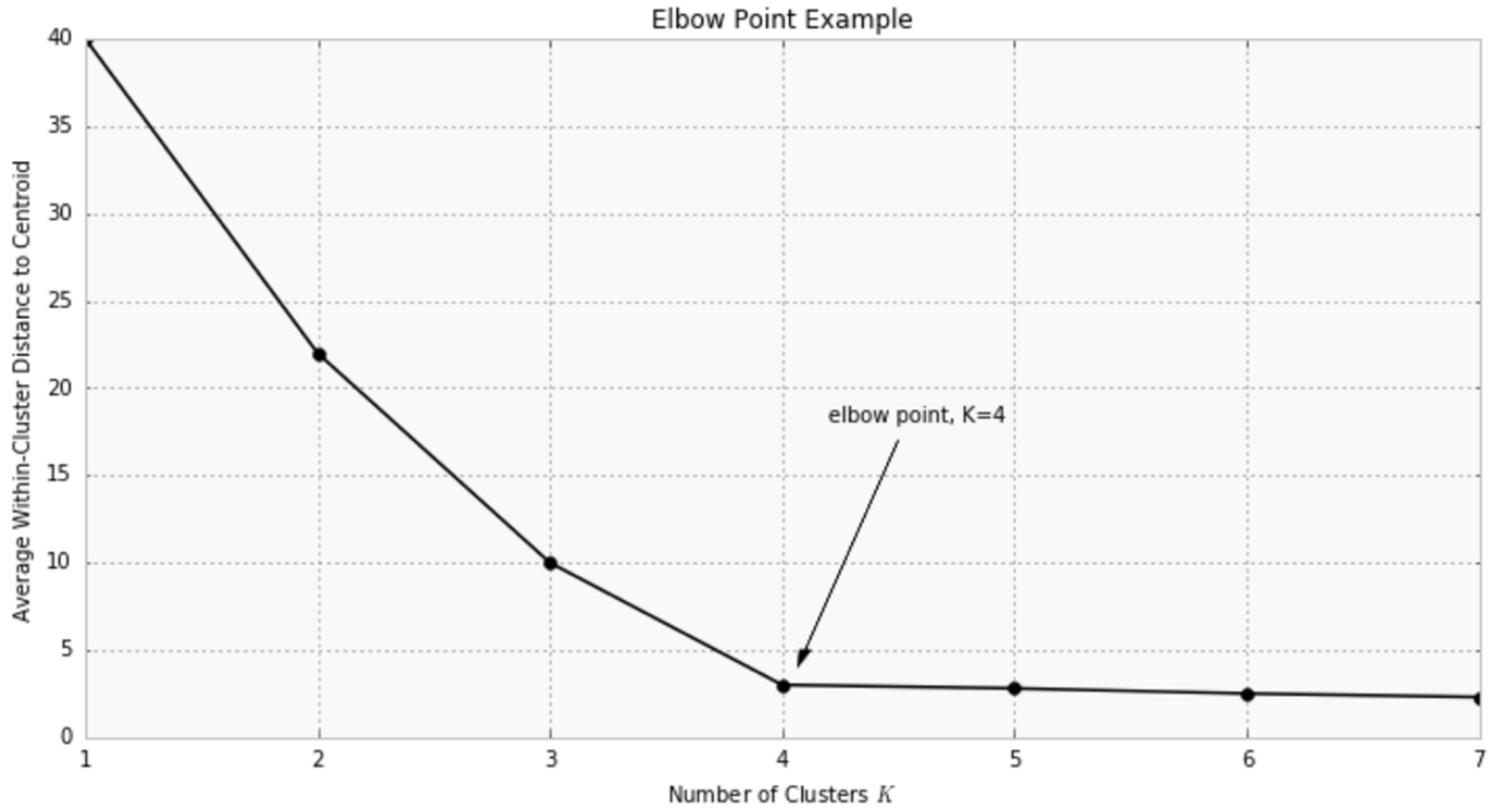
AGRUPAMIENTO JERARQUICO

# DENDOGRAMA



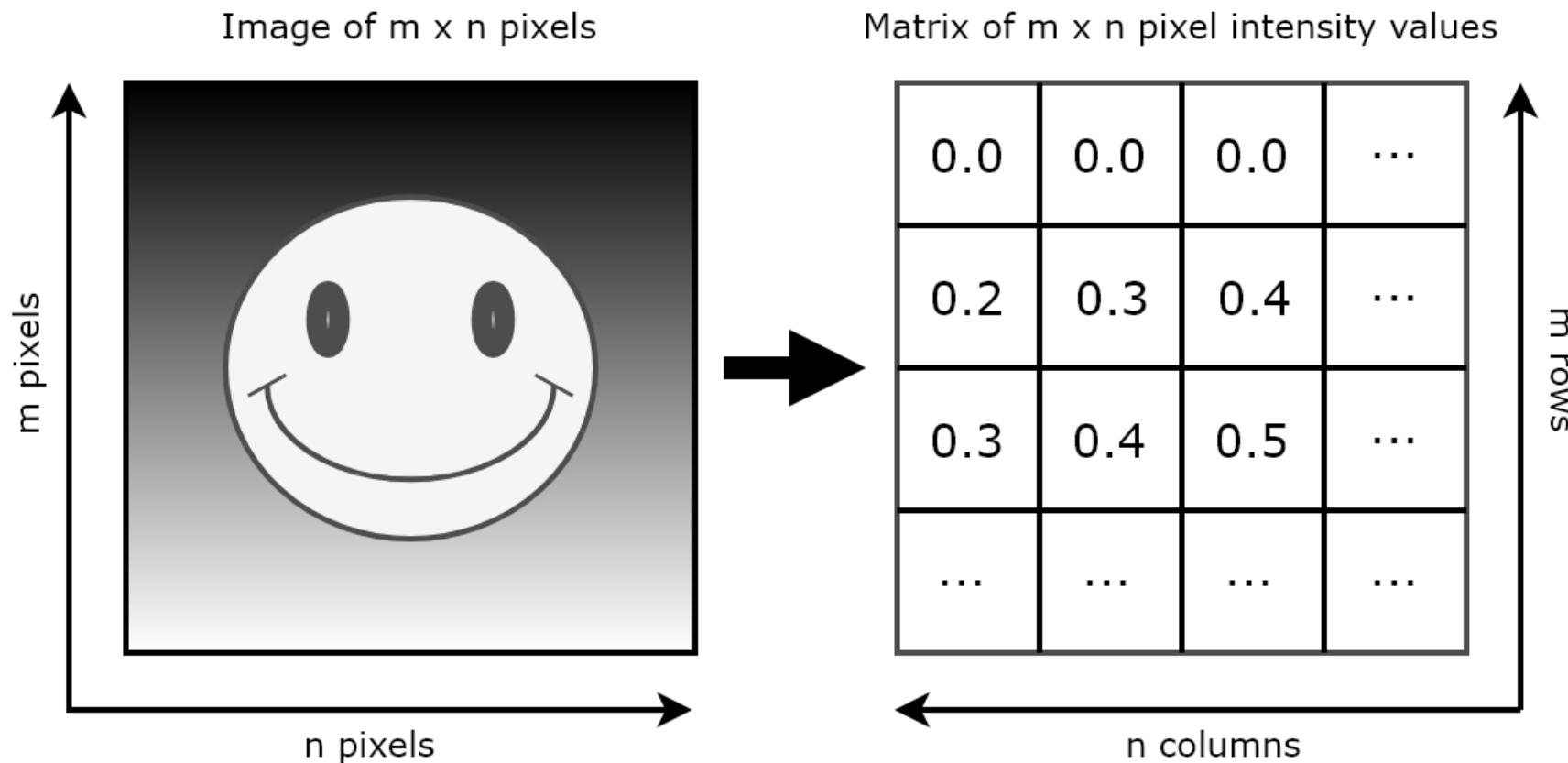
# K-MEDIAS





# ELBOW METHOD

# IMAGEN = MATRIZ



```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK

# SISTEMAS DE RECOMENDACIÓN





What to Watch

My Subscriptions

Music



[Last Week Tonight with John Oliver: Online Harassment...](#)  
by LastWeekTonight  
2,895,113 views • 1 week ago

[BASTILLE feat. Ella - No Angels \(HD\)](#)  
by dancenationedm  
5,836,404 views • 2 years ago

[Bill Maher with Charlie Rose](#)  
by Darklordabc  
88,923 views • 8 months ago



[START! Walking at Home](#)  
American Heart Association  
Leslie Sansone's Walk at Home  
63,839 views • 7 months ago

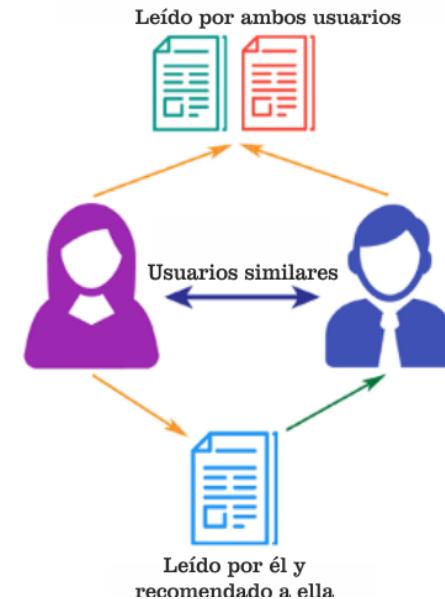
[The 5 Best Shia Labeouf "DO IT" Motivation Videos](#)  
by The VGC  
1,390,971 views • 4 weeks ago

[Elon Musk and Bill Gates discuss AI, entrepreneurship...](#)  
by Every Elon Musk Video  
105,539 views • 2 months ago

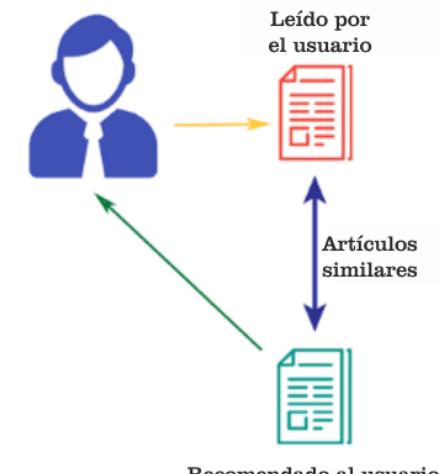
# CASOS DE USO

# TIPOS DE SISTEMAS DE RECOMENDACIÓN

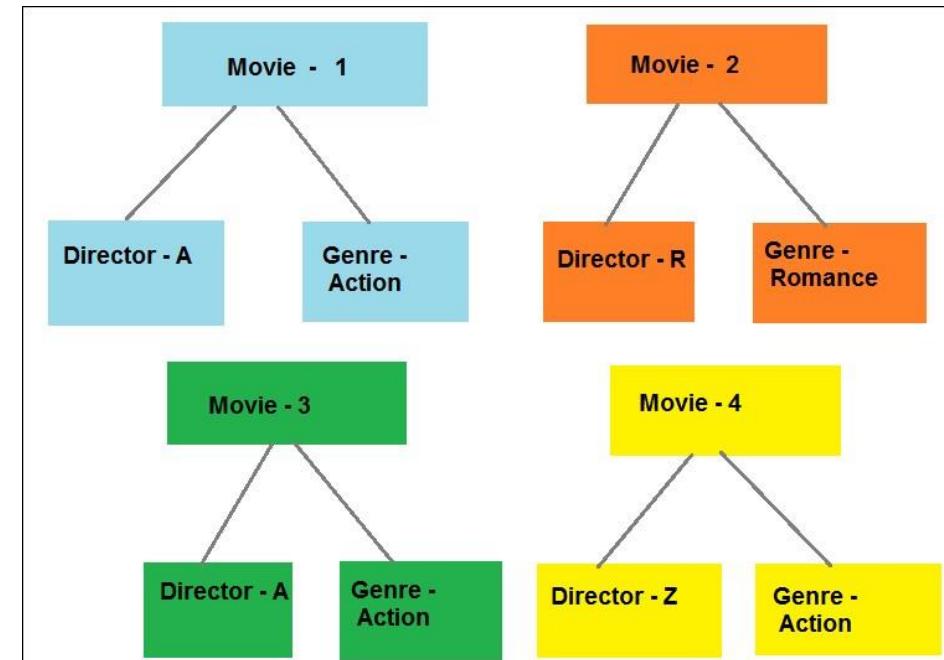
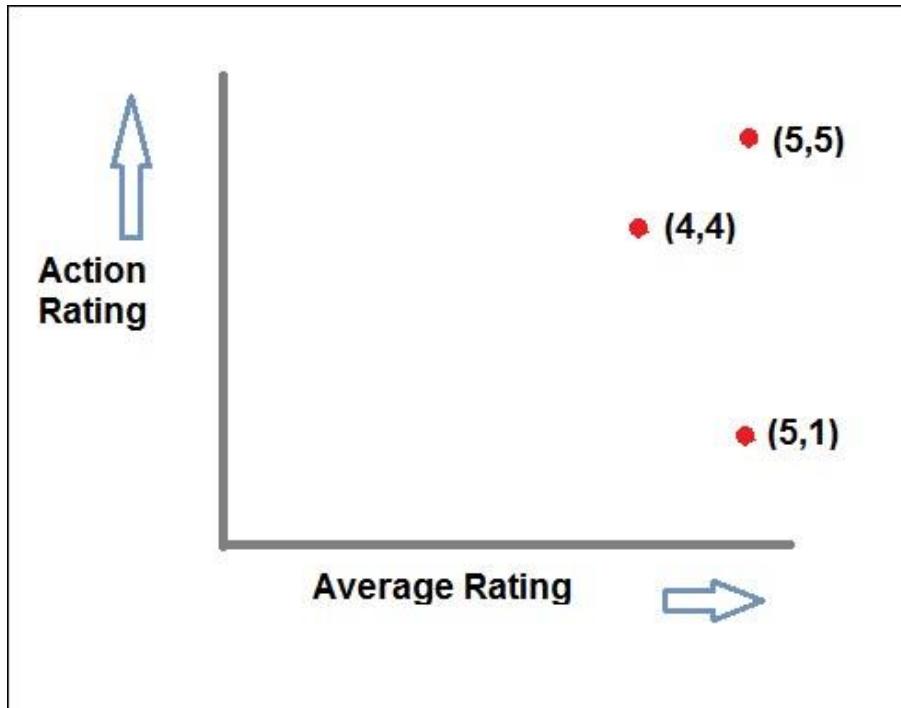
**Filtro colaborativo**



**Filtro basado en contenido**



# BASADOS EN CONTENIDO





Similar Users, User-A and User-B with the movies they have watched

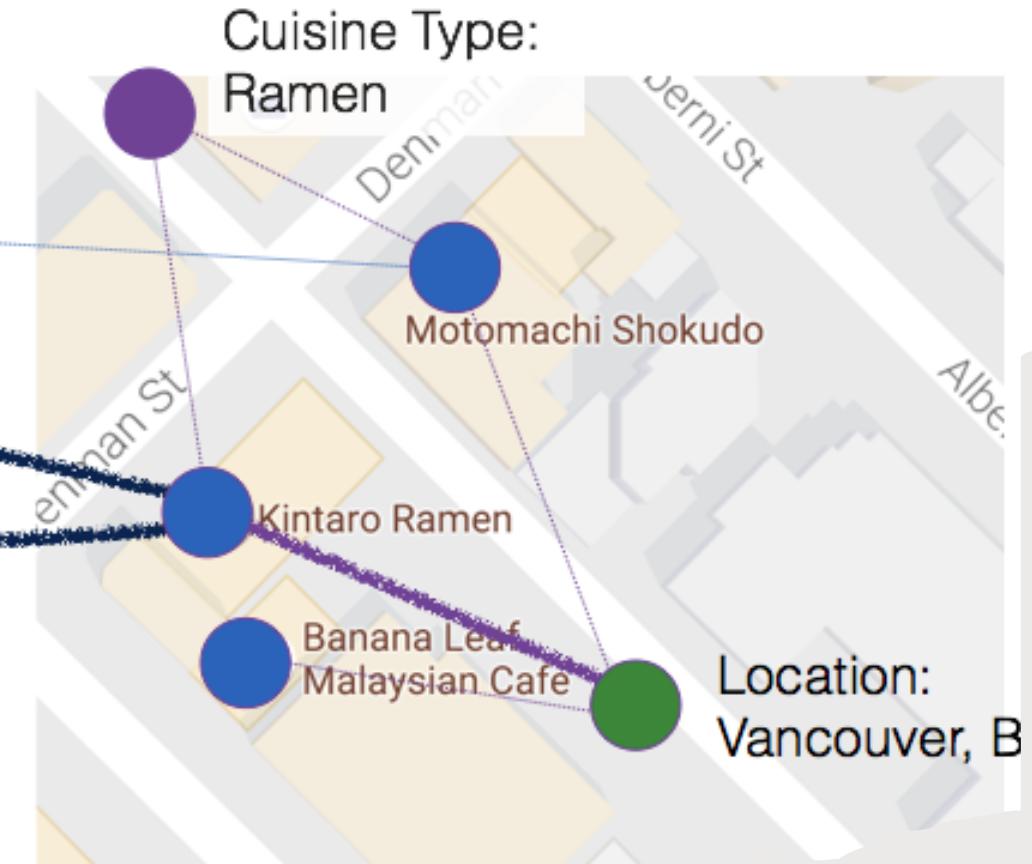
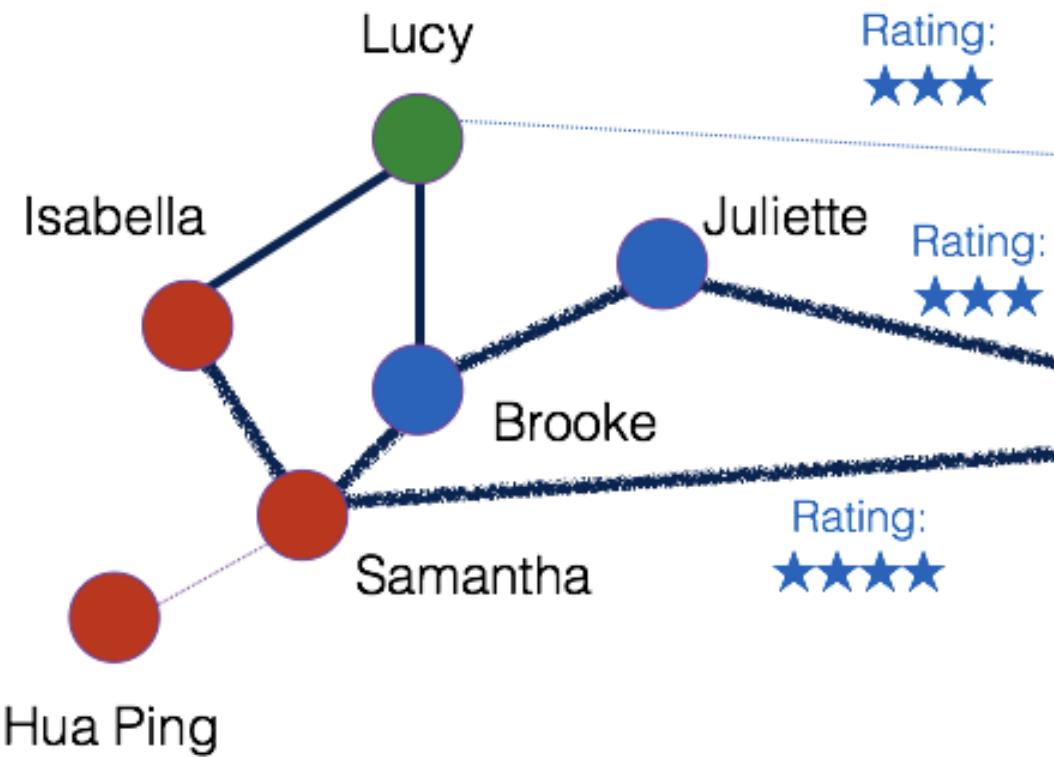
# FILTRO COLABORATIVO

```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

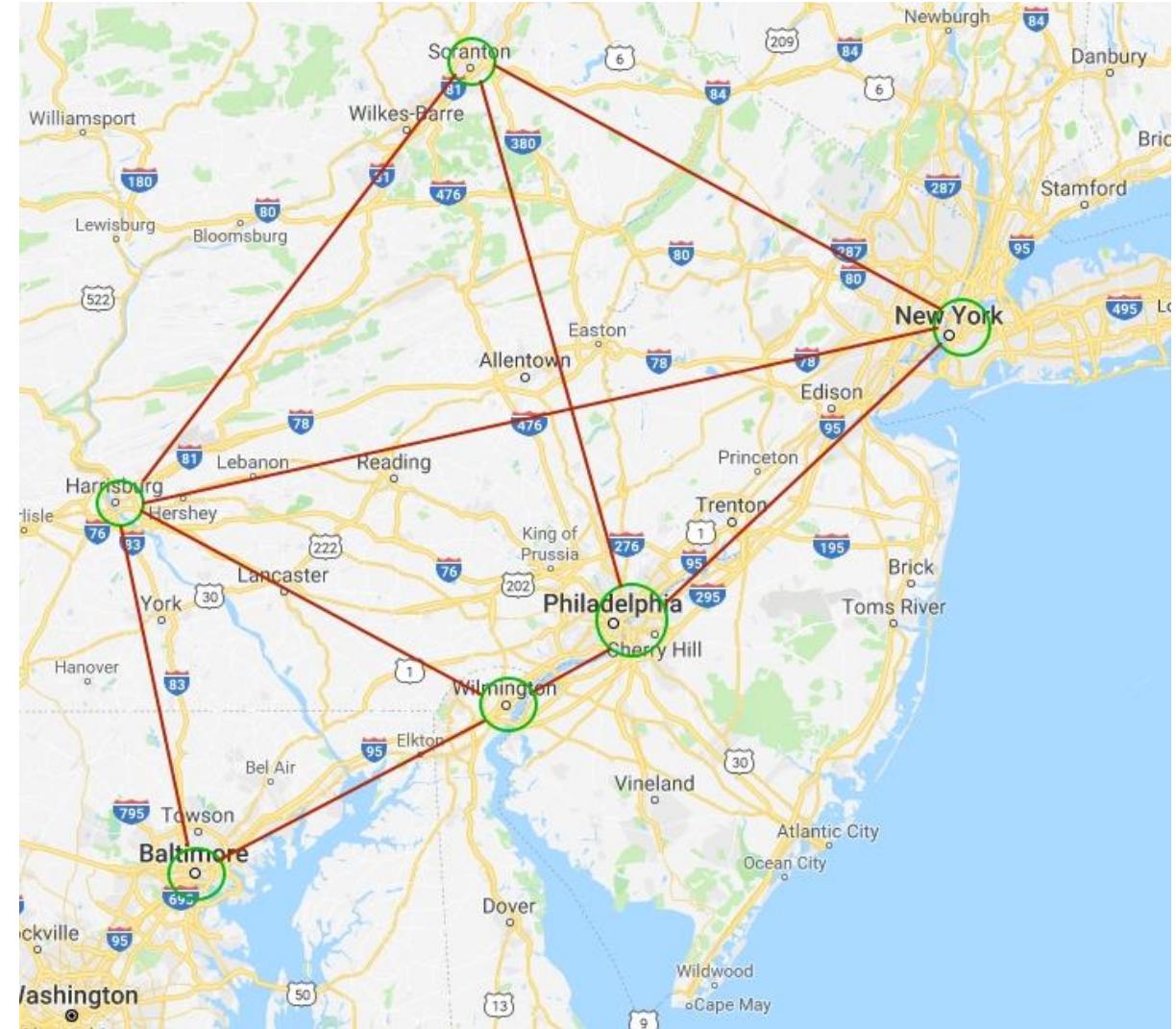
# NOTEBOOK



## social network + restaurant recommendations

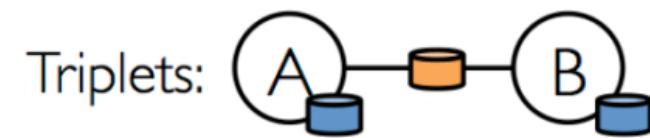


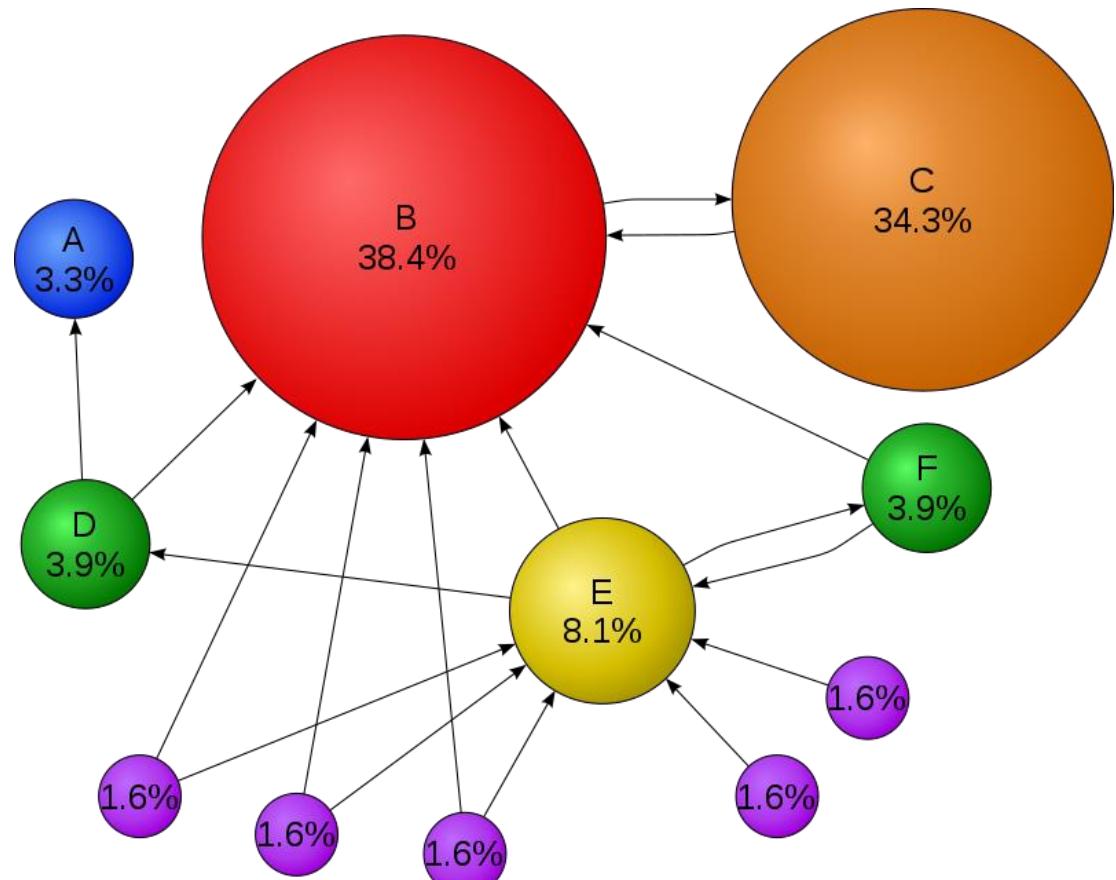
# GRAFO DEL EJEMPLO



```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54         });
55       }
56     });
57   });
58 }
```

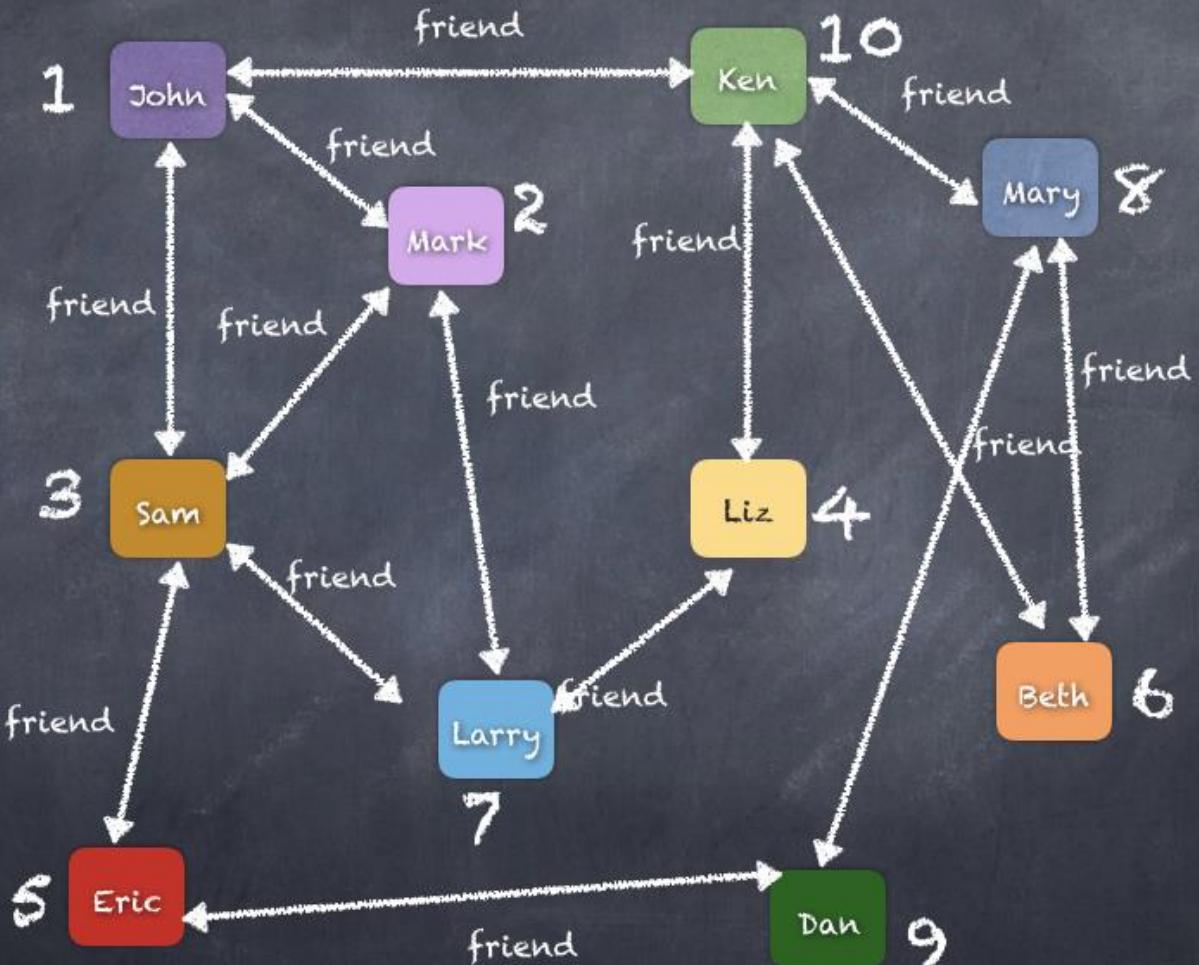
## LABORATORIO 11



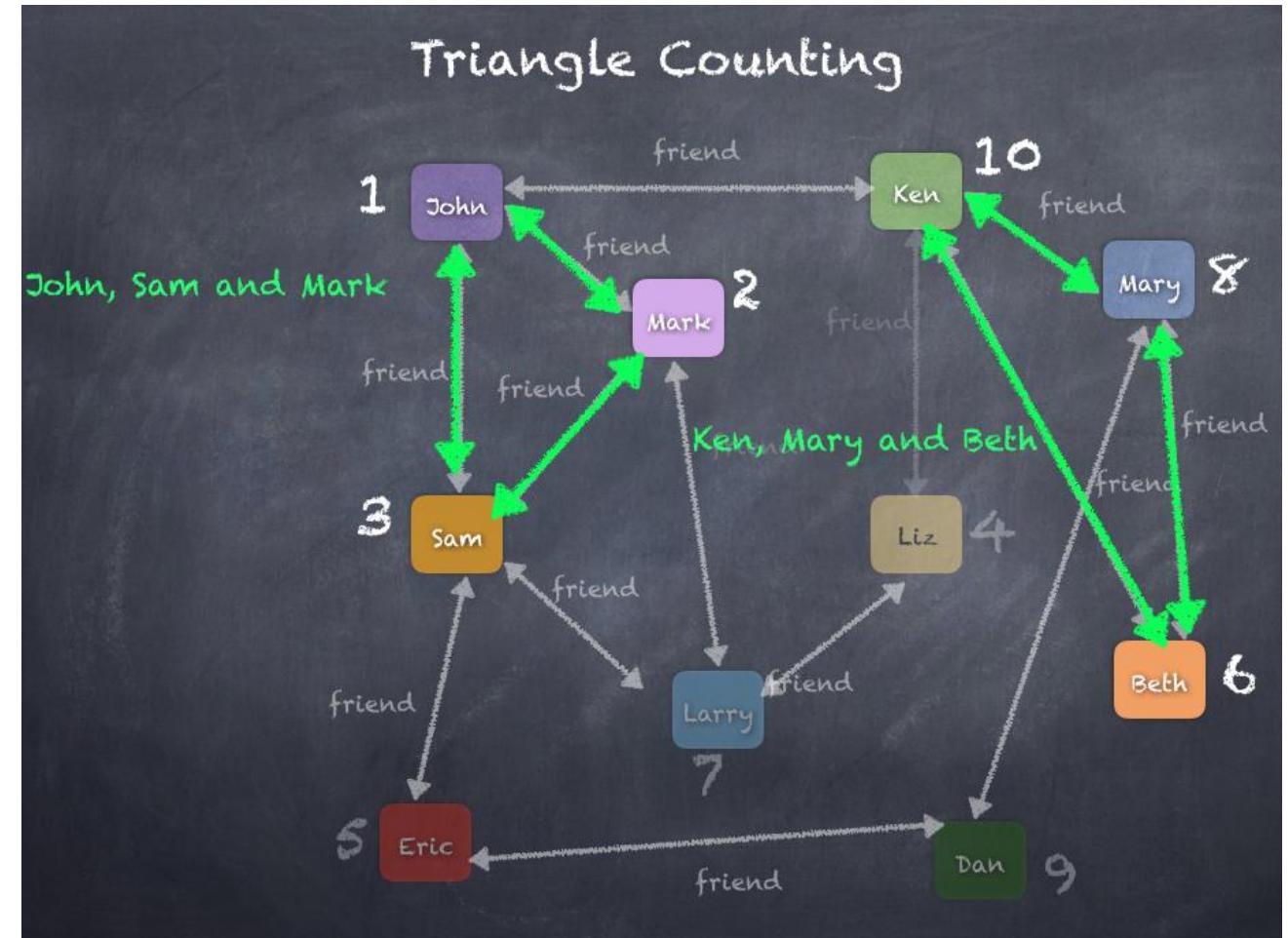


# ALGORITMOS

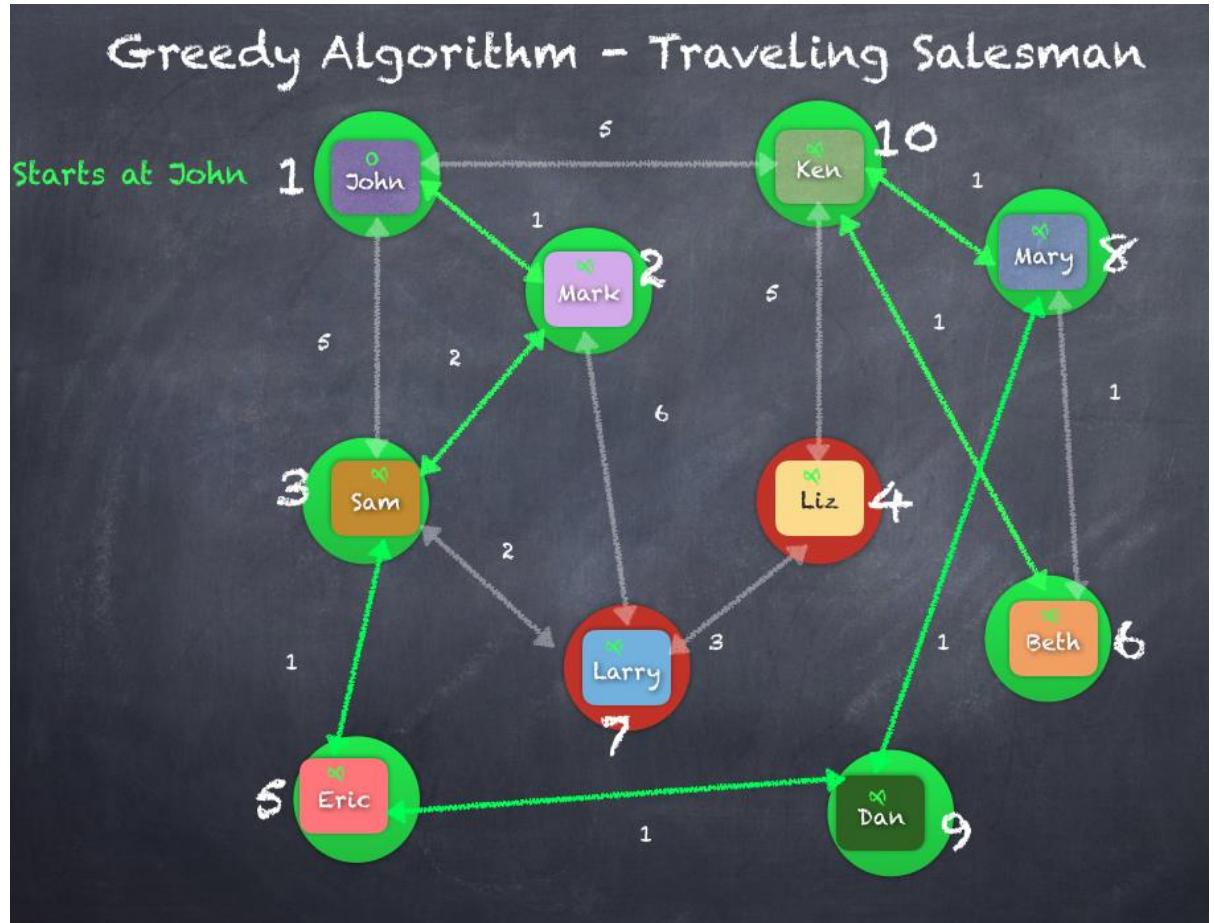
## Graph of Facebook Friends



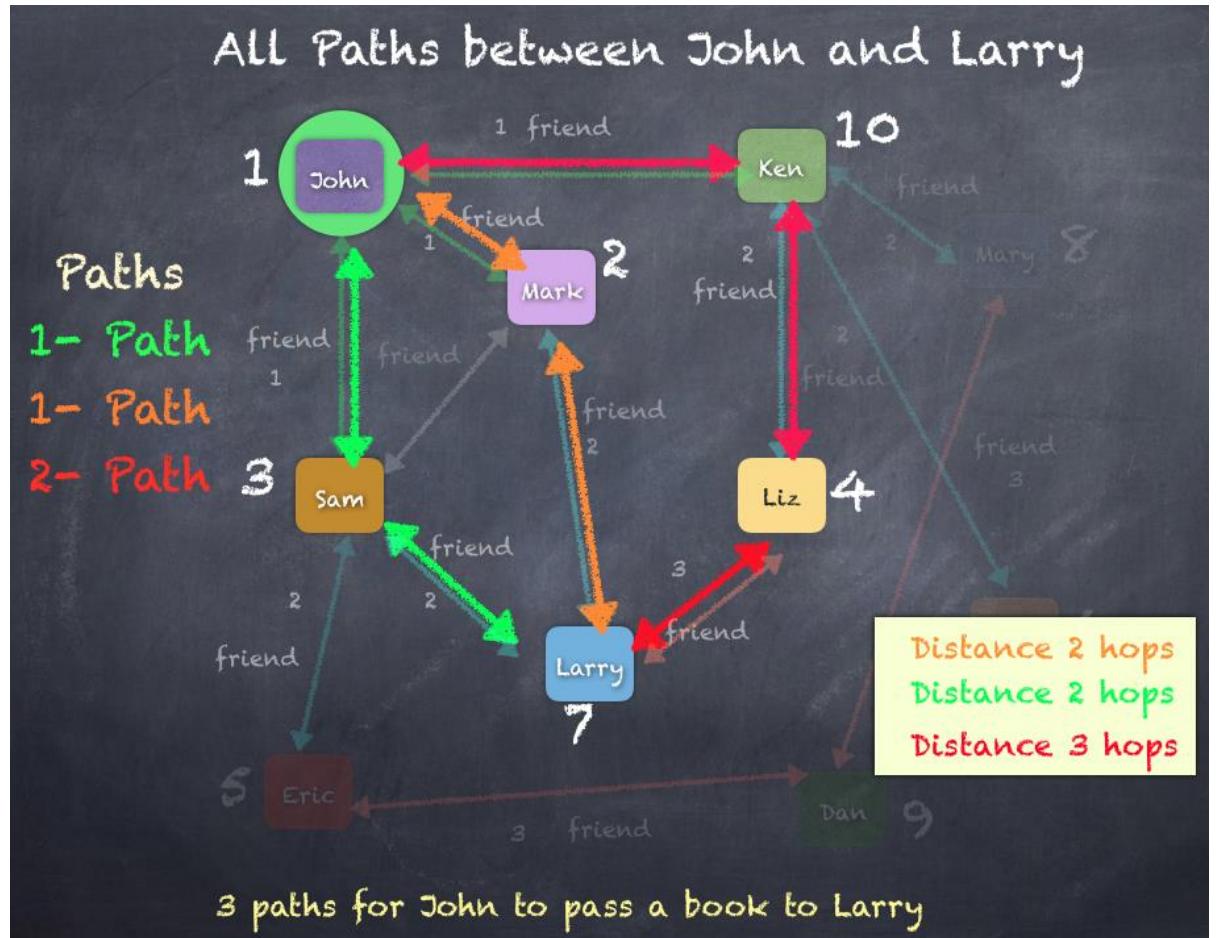
# TRIANGULOS



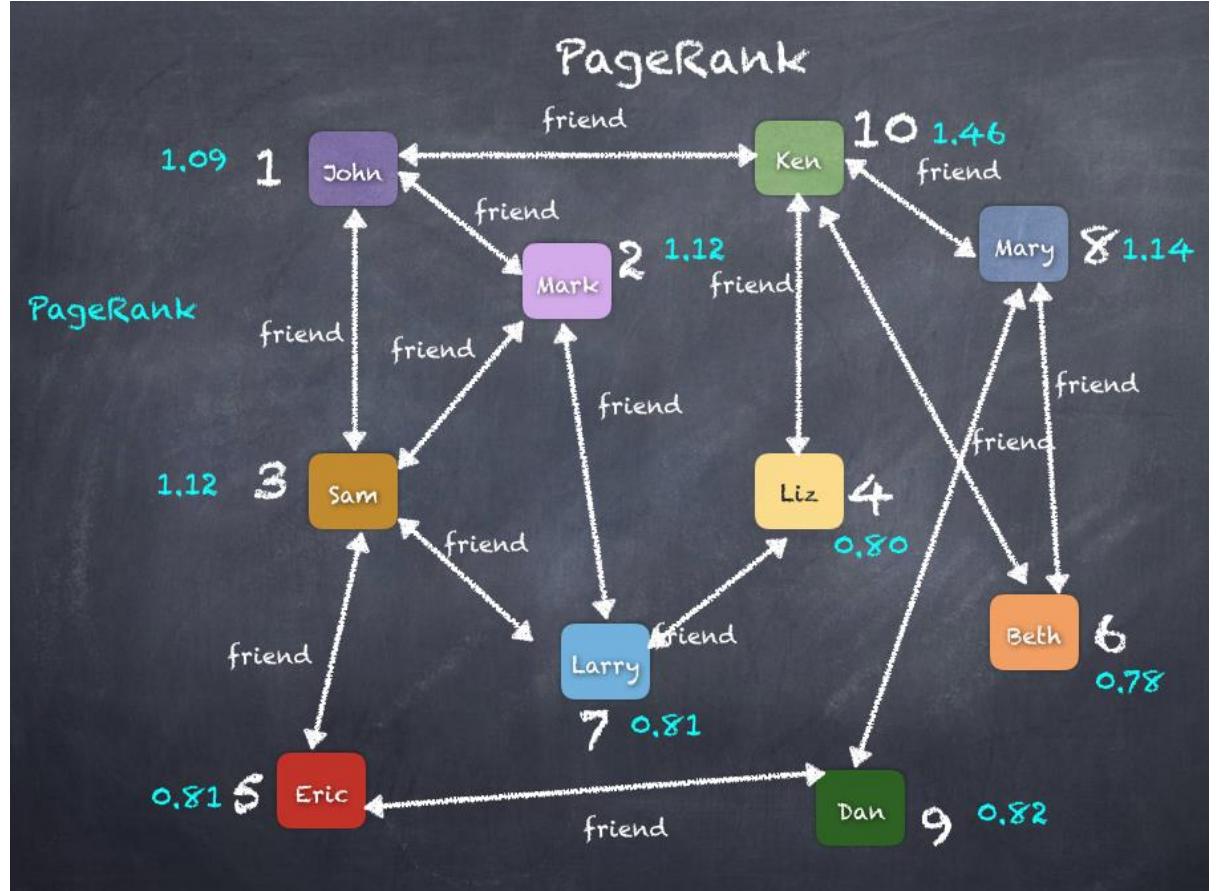
# PROBLEMA DEL VENDEDOR AMBULANTE



# RUTA MAS CORTA



# PAGE RANK

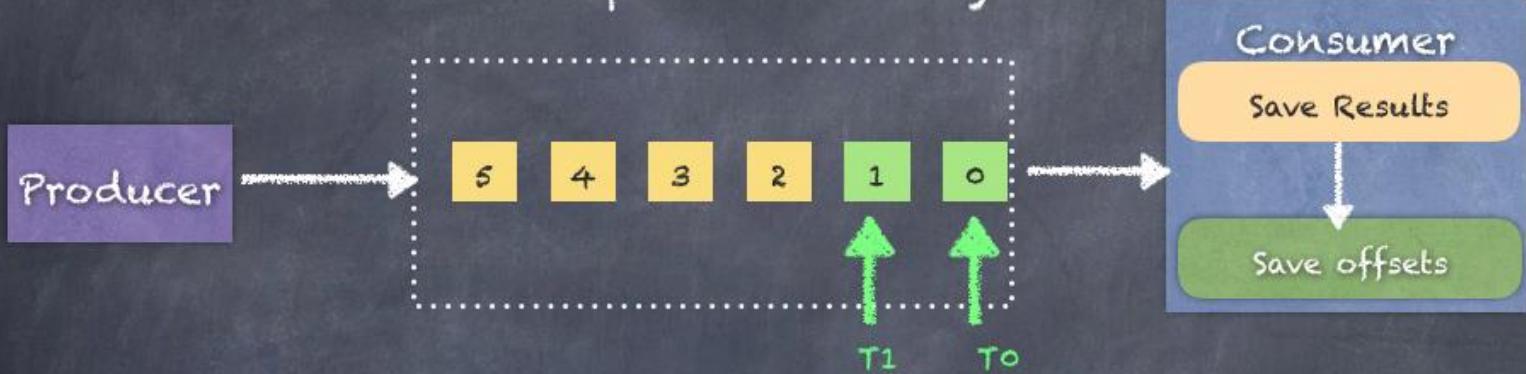


# STREAMING



# At Least Once Processing

events are processed one by one

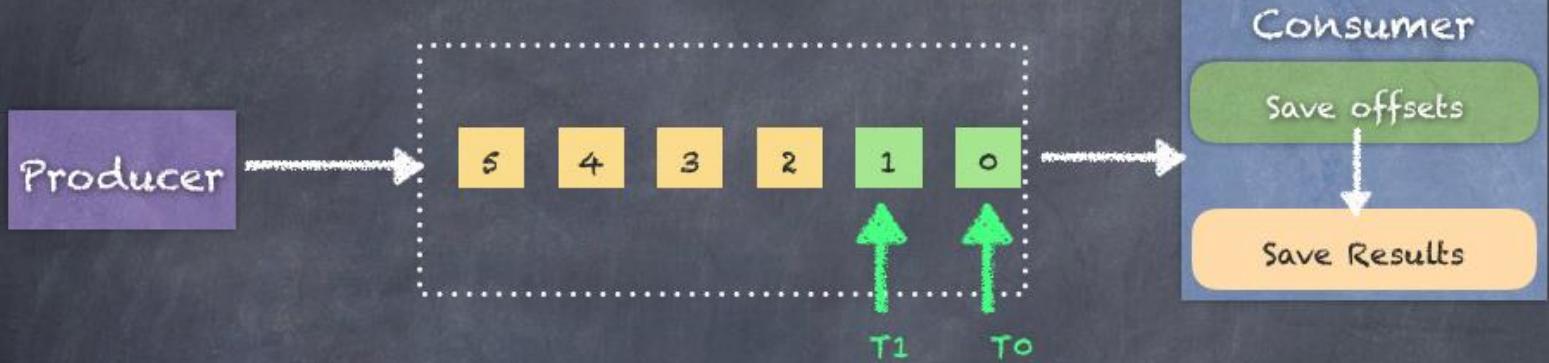


On Failure, event 0 is processed twice

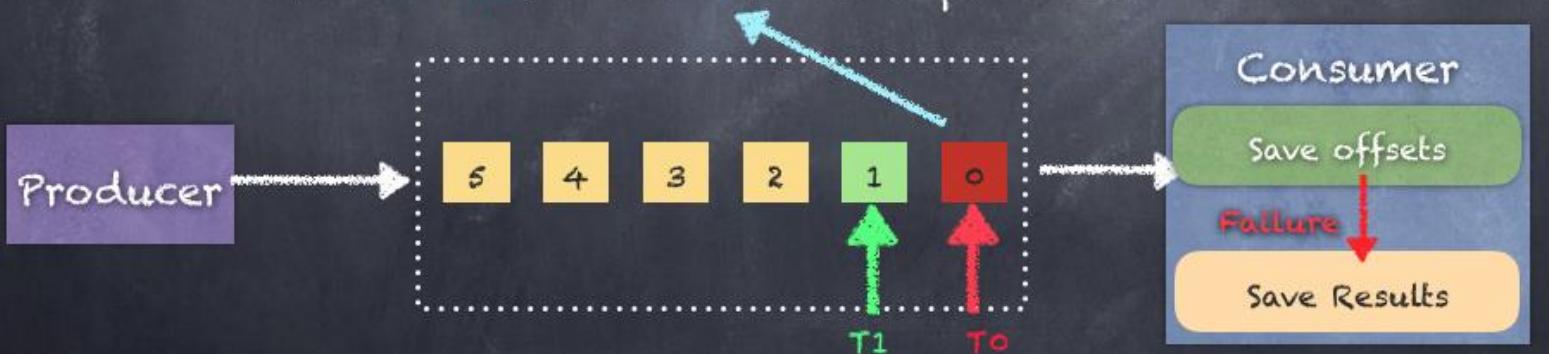


# At Most Once Processing

events are processed one by one

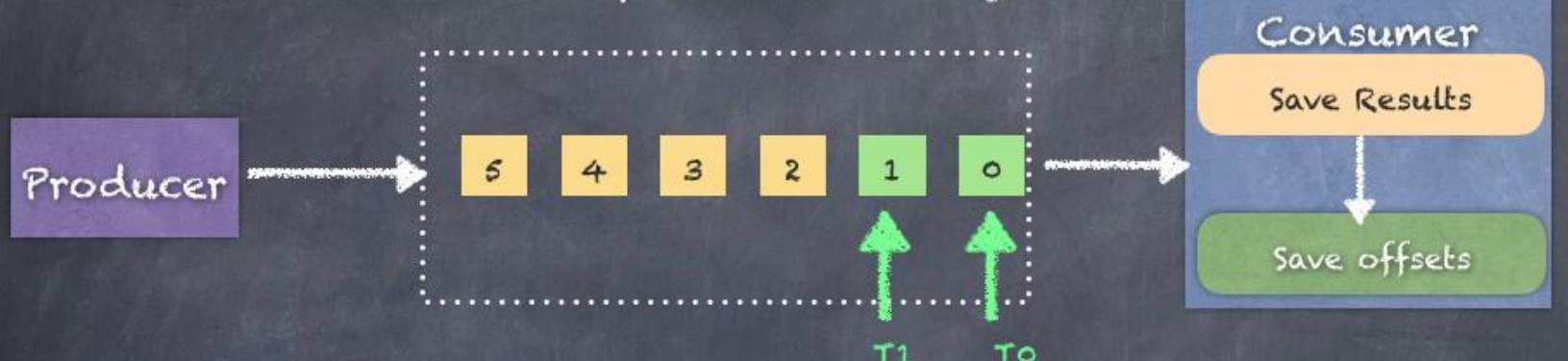


On Failure, event 0 is never processed

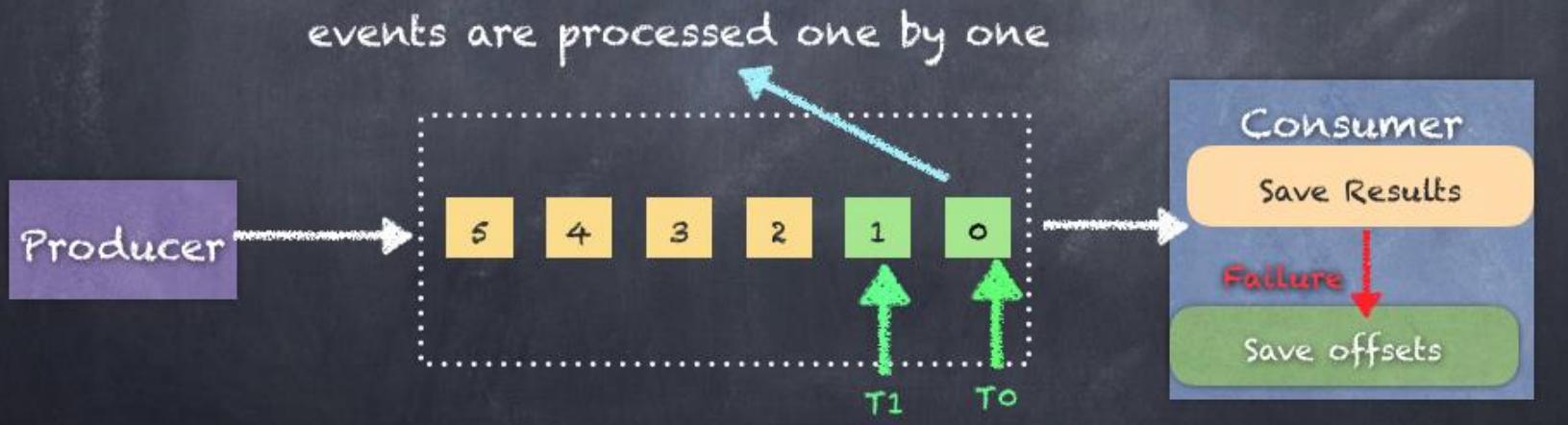


# Exactly Once Processing

events are processed one by one



events are processed one by one

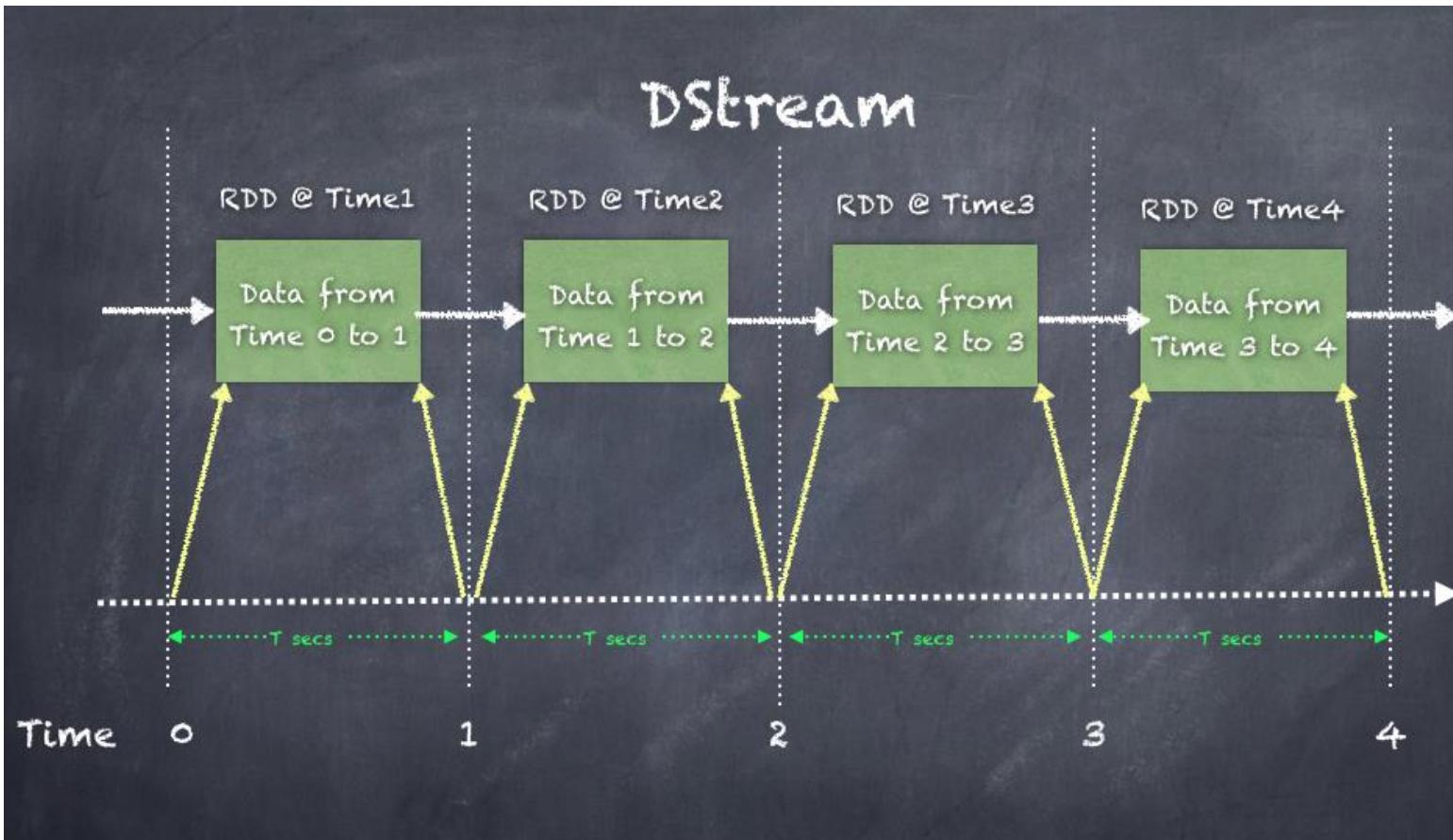




SPARK STREAMING

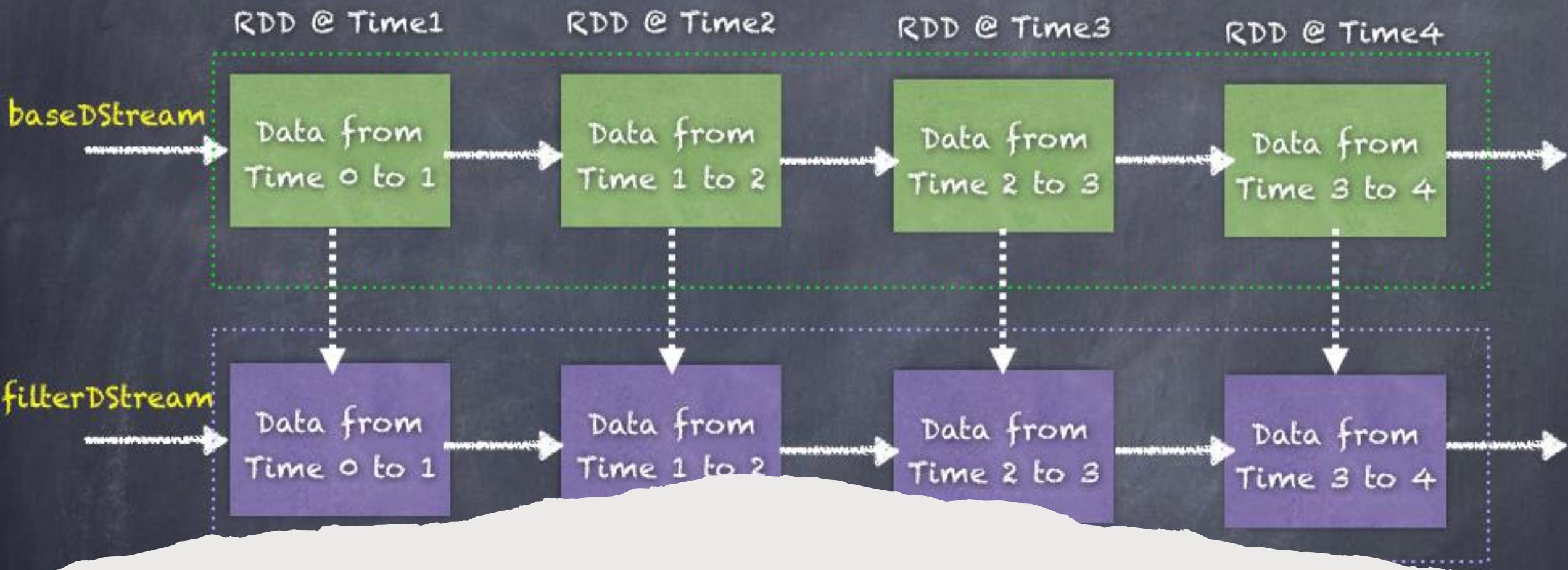


MICROOTES



# DSTREAMS

# DStream transformation



TRANSFORMACIONES

```
36
37     'drag',
38     'drop',
39   ];
40   events.forEach(e => {
41     fileDropZone.addEventListener(e, (ev) => {
42       ev.preventDefault();
43       if (ev.type === 'dragenter') {
44         fileDropZone.classList.add('solid-border');
45       }
46       if (ev.type === 'dragleave') {
47         fileDropZone.classList.remove('solid-border');
48       }
49       if(ev.type === 'drop') {
50         fileDropZone.classList.remove('solid-border');
51         ev.dataTransfer.files
52         .values.map(tag => {
53           tag.setAttribute('class', 'tag');
54           tag.setAttribute('border', '1px solid black');
55           tag.setAttribute('border-radius', '10px');
56         })
57       }
58     });
59   });
60 }
```

# NOTEBOOK